

Univariate, bivariate, and multivariate methods in corpus-based lexicography – a study of synonymy

Antti Arppe

*Academic dissertation to be publicly discussed, by due permission of the Faculty of
Arts at the University of Helsinki in lecture room 13, on the 19th of December, 2008,
at 12 o'clock.*

University of Helsinki
Department of General Linguistics
P.O. Box 9 (Siltavuorenpenger 20 A)
FI-00014 University of Helsinki
Finland

PUBLICATIONS
NO. 44
2008

Cover image:
The tree seen from the researcher's chamber on 6.5.2008 (Antti Arppe)

ISSN 0355-7170
ISBN 978-952-10-5174-6 (paperback)
ISBN 978-952-10-5175-3 (PDF)

URL: <http://ethesis.helsinki.fi/>

Helsinki 2008
Helsinki University Print

*To my father and mother
Juhani and Raija Arppe*

Abstract

In this dissertation, I present an overall methodological framework for studying linguistic alternations, focusing specifically on lexical variation in denoting a single meaning, that is, synonymy. As the practical example, I employ the synonymous set of the four most common Finnish verbs denoting THINK, namely *ajatella*, *mieltiä*, *pohtia* and *harkita* ‘think, reflect, ponder, consider’. As a continuation to previous work, I describe in considerable detail the extension of statistical methods from dichotomous linguistic settings (e.g., Gries 2003; Bresnan et al. 2007) to polytomous ones, that is, concerning more than two possible alternative outcomes.

The applied statistical methods are arranged into a succession of stages with increasing complexity, proceeding from univariate via bivariate to multivariate techniques in the end. As the central multivariate method, I argue for the use of polytomous logistic regression and demonstrate its practical implementation to the studied phenomenon, thus extending the work by Bresnan et al. (2007), who applied simple (binary) logistic regression to a dichotomous structural alternation in English.

The results of the various statistical analyses confirm that a wide range of contextual features across different categories are indeed associated with the use and selection of the selected think lexemes; however, a substantial part of these features are not exemplified in current Finnish lexicographical descriptions. The multivariate analysis results indicate that the semantic classifications of syntactic argument types are on the average the most distinctive feature category, followed by overall semantic characterizations of the verb chains, and then syntactic argument types alone, with morphological features pertaining to the verb chain and extra-linguistic features relegated to the last position.

In terms of overall performance of the multivariate analysis and modeling, the prediction accuracy seems to reach a ceiling at a *Recall* rate of roughly two-thirds of the sentences in the research corpus. The analysis of these results suggests a limit to what can be explained and determined within the immediate sentential context and applying the conventional descriptive and analytical apparatus based on currently available linguistic theories and models.

The results also support Bresnan’s (2007) and others’ (e.g., Bod et al. 2003) probabilistic view of the relationship between linguistic usage and the underlying linguistic system, in which only a minority of linguistic choices are categorical, given the known context – represented as a feature cluster – that can be analytically grasped and identified. Instead, most contexts exhibit degrees of variation as to their outcomes, resulting in proportionate choices over longer stretches of usage in texts or speech.

Preface and acknowledgments



ud mu-e-ši-zal

a-na-am₃ šu mu-da-ti

‘The time passed, and what did you gain?’^{1 2}



nam-sag₉-ga kaš-a

nam-ḥul kaskal-la

‘The good thing – it is the beer,
the bad thing – it is the journey.’³

This dissertation is essentially the product of a decade-long sequence of innumerable discussions, in which my supervisors, colleagues, friends, and family have imparted ideas and insights that have edged my research further. Not only have these discussions taken place during supervision and consultation sessions proper, but they have chanced to happen while crossing a street, in informal meetings in the University’s corridors, by the copy-machine or printer at the Department, over the phone or via e-mail, during the short question-and-answer sessions or the longer chit-chat that follows at academic conferences, workshops, and seminars, or even through the age-old medium for the exchange of scientific ideas through written books or articles.

All too often, I have not fully appreciated such ideas when they have first been proffered to me – probably because I was not yet ready to understand them yet – and sometimes the true value of many a suggestion has dawned on me only years afterwards. In all likelihood, my partners in these discussions may not have realized themselves what – in retrospect – crystal-clear one-liners they have uttered or written concerning one aspect or another relevant to my research process. Nonetheless, these giveaway reflections are what I can now see to form the backbone of this written work. Although this study is thus fundamentally the result of applying an ensemble of many borrowed ideas, my own contribution – and therefore also nobody else’s responsibility but mine – is the complex whole that they constitute and the overall interpretation that is conveyed.

¹ *Electronic Text Corpus of Sumerian Literature: Proverbs, Collection 3: 3.157* (Black et al. 1998-2008).

² *CompositeCuneiform signs by Tinney and Everson (2004-2007)*.

³ *Electronic Text Corpus of Sumerian Literature: Proverbs, Collection 7: 7.98* (Black et al. 1998-2006).

I am deeply grateful to my four official supervisors, Fred Karlsson, Lauri Carlson, Urho Määttä, and Juhani Järvikivi, for each guiding in their own characteristic manner this work further towards its final conclusion. Fred Karlsson saw promise in my early, quite sketchy research plan concerning synonymy, accepting me, as the academic immigrant from engineering that I was, into the field of linguistics as a post-graduate student, as well as tipping me off about my first research funding in the *GILTA* project. Fred's persistent insistence on focusing on the final goal, however uncomfortable I found the issue at times, was instrumental in getting this dissertation finished.

Lauri Carlson was a wizard in articulating what I was not yet able to, seeing where my thinking was taking me before I knew it myself, and crystallizing these into simple comments or questions that connected my tiny, seemingly isolated island of research into the grander scheme of things within the linguistic discipline. During the final writing process of this text, Lauri pored through diligently each chunk as I succeeded in churning them out, providing me the assurance that I remained on the right track.

Urho Määttä always had ample time and interest to discuss the most general ramifications and connections of my work and to delve deep into what language as an object of study and linguistics as a discipline were fundamentally about. At the same time, Urho assisted me in contenting myself first with only one, seemingly small, part of my original research plan, and then mastering that portion as comprehensively as possible.

Juhani Järvikivi's expertise in psycholinguistic experimentation was the component that I did not know was missing until I stumbled onto it, and which was to open my eyes to what my research in linguistics could on the long term concern. Co-authoring two journal articles together with Juhani introduced me to what academic discourse at the international level is really about, and this co-operation also allowed me to sketch out the general, multimethodological backdrop which this considerably more tightly focused work serves.

I am also indebted to Martti Vainio, who acted as my fifth supervisor in all but name. It was Martti who realized that logistic regression might be the very statistical method that could bring some order to my jungle of linguistic features, and he neither hesitated nor spared his time in working out with me a solution to applying this technique to my multiple-outcome setting and getting me started with *R*. I also appreciate Martti's iconoclastic attitude in the many intensive but free-wheeling discussions we have had concerning a wide array of topics ranging from the current central questions in linguistics and other sciences to international politics, in which we have not been troubled by the passing of time.

Furthermore, I am grateful to Kimmo Koskenniemi for all the support and sincere attention he has shown to the progress of my research work. Thanks to Kimmo, I originally wound up working in the mid-to-late 1990s at Lingsoft, a small Finnish language technology company, where one of the software development projects I was involved in, namely *inflecting thesauri*, would lead me to discover the kernel of this dissertation. In particular, I want to thank Kimmo for inviting me on several occasions

to speak in the language technology seminar in order to sort out the state of my research and rediscover its red thread, when I was in danger of losing my way.

At the Department of General Linguistics, I have appreciated the collegial and informal atmosphere, which has allowed me to benefit from the experience of the research and teaching staff representing three closely related but distinct subjects. I am thankful for the patient coaching and support by Jan-Ola Östman and Kari K. Pitkänen when I was still primarily a full-time novice in the field. Moreover, I have fond memories of the late Orvokki Heinämäki, who taught some of the very first courses in general linguistics that I took in the early 1990s. Most importantly, I have also had the excellent opportunity to broaden my understanding of linguistics through encountering the extraordinarily diverse kaleidoscope research topics pursued by the Department's post-graduate students and affiliated post-doctoral researchers. Among many others, I am happy to have had a lively interchange with Matti Miestamo, in which he has convincingly argued for the importance of typological perspective in linguistics. In particular, I want to thank Matti for providing me with comments in this respect on the Introduction and Discussion of this text, and for the extra-curricular discussions we have had as fellow Kallio linguists.

I am very grateful to both of my preliminary examiners, R. Harald Baayen (University of Alberta) and Stefan Th. Gries (University of California, Santa Barbara), for their constructive criticism with respect to the current text and their encouraging suggestions to continue developing its themes further. In fact, when I had first made the acquaintance of both my examiners, on separate occasions just around the beginning of this millennium, they had each given me tips that I did not heed at the time, but which advice I later realized to contain two cornerstones of this dissertation. Harald Baayen was the first person to recommend to me that I get acquainted with the *R* statistical programming environment, which became the workhorse for all the statistical analysis in my work, while Stefan Th. Gries' own doctoral dissertation already contained the general tripartite methodological setup that I was to find well-suited to bring structure also to my work.

In addition, there are a great number of scholars both in Finland and abroad with whom I have had the opportunity, privilege and pleasure of collaborating or exchanging ideas, opinions, and assistance over the many years, not always only concerning linguistics but also the practical everyday challenges of leading the life a young researcher. These people whom I wish to recognize are (in alphabetical order): Daniel Aalto, Tiina Arppe, Lili Aunimo, Juhani Birn, Pia Brandt, Joan Bresnan, Andrew Chesterman, Pernilla Danielsson, Dagmar Divjak, Antonina Durfee, Nils-Erik Enkvist, Stefan Evert, Marja Etelämäki, Sam Featherston, Dylan Glynn, Stefan Grondelaers, Mickel Grönroos, Erja Hannula, Tarja Riitta Heinonen, Irmeli Helin, Kris Heylen, Suvi Honkanen, Timo Honkela, Silja Huttunen, Esa Itkonen, Jarmo H. Jantunen, Eriika Johansson, Kristiina Jokinen, Timo Järvinen, Panu Kalliokoski, Irina Kauhanen, Tapani Kelomäki, Harri Kettunen, Tarja Knuutila, Leena Kolehmainen, Lari Kotilainen, Ville Laakso, Krista Lagus, Ritva Laury, Jaakko Leino, Pentti Leino, Yrjö Leino, Krister Lindén, Mikko Lounela, Aapo Länsiluoto, Laura Löfberg, Anke Lüdeling, Annu Marttila, Matti Miestamo, Manne Miettinen, Sjur Nørstebø Moshagen, Jussi Niemi, Otto Nieminen, Urpo Nikanne, Alexandre Nikolaev, Torbjørn Nordgård, Elina Nurmi, Martti Nyman, Seppo Nyrkkö, Krista Ojutkangas, Jussi Pakkasvirta, Heikki Patomäki, Pertti Palo, Santeri Palviainen, Simo Parpola,

Marja Peltomaa, Kaarina Pitkänen, Marja Pälsi, Michaela Pörn, Jarno Raukko, Timo Riiho, Jouni Rostila, Jack Rueter, Janne Saarikivi, Gabriel Sandu, Dirk Speelman, Kaius Sinnemäki, Pirkko Suihkonen, Mickael Suominen, Antti Suni, Gert De Sutter, Pasi Tapanainen, Lauri Tarkkonen, Jarmo Toivonen, Trond Trosterud, Jarno Tuimala, José Tummers, Tuuli Tuominen, Ulla Vanhatalo, Johanna Vaattovaara, Kari T. Vasko, Kimmo Vehkalahti, Marja Vierros, Simo Vihjanen, Liisa Vilkki, Maria Vilkuna, Laura Visapää, Mari Voipio, Atro Voutilainen, Tytti Voutilainen, Jarmo Välikangas, Fredrik Westerlund, Caroline Willners, Malene Würtz, Anssi Yli-Jyrä, and Jussi Ylikoski.

From my time at Lingsoft, I am especially grateful to Juhani Birn for his witty insight and seasoned experience in interpreting the unexpected behavior of inflecting thesauri. Likewise, I appreciate the initial effort that Mari Voipio and Malene Würtz made to put our finger on what the oddities exactly were and what they might generally entail. Moreover, I truly enjoyed the many manifold times that we managed to drift over to discussing linguistic topics together with many of the company's employees in the late 1990s.

For the first three years of my postgraduate studies, I participated in and received funding through the *GILTA* project, led by Ari Visa (Tampere University of Technology), Hannu Vanharanta (Tampere University of Technology), and Barbro Back (Åbo Akademi University). I appreciate greatly the wide latitude that I was allowed within this cross-disciplinary project, permitting me to pursue my purely linguistic research interests at the same time. Especially glad I am that we succeeded in tying most of the different strands of research together in one major joint article, in which the individual contributions of Tomas Eklund (Åbo Akademi University) and Camilla Magnusson played a central role. Although Camilla started out as my research assistant in the *GILTA* project, she soon became an independent-minded researcher in her own right, and I am grateful to her for questioning some certainties that I held.

Already in 1999, I had been invited by Anu Airola and Hanna Westerlund to join an informal post-graduate discussion group, which soon extended to include also Jussi Piitulainen. I am thankful to Anu, Hanna, and Jussi for all the mutual assistance and instruction that was of great help to me in trying to make sense – both to myself and the others – of what various central theoretical concepts meant or could connote. During 2002-2003, this evolved into the *Suomenlinna* circle of linguists, also known as *TR-9*, formed by Urho Määttä, Anu Airola, Juhani Järvikivi, Camilla Magnusson, Jussi Piitulainen, Reetta Konstenius, Martti Vainio, Hanna Westerlund, and myself. *TR-9* gave me an inspiring glimpse of the satisfying intellectual interchange that can be achieved when linguists with diverse backgrounds and complementary skills work together rather than alone. In particular, I still remember clearly Reetta's insistence on the rigorous formation of hypotheses and firm argumentation while conducting linguistic research.

Beginning in 2003, I had the privilege of being accepted into *LANGNET*, the Finnish Graduate School in Language Studies, within which I received full funding for the first four years, followed by an additional year as a self-financed fellow. I have come to consider many of my postgraduate colleagues within the *LANGNET* graduate school, as well as recent doctors, to form a unique peer group that is best positioned to

wholly understand the many conflicting demands and to provide support in the difficult challenges that young researchers face these days. Many of these originally primarily professional colleagues I now consider also my friends.

During my research process and the writing of this dissertation text, I have had the privilege of tapping into the expert knowledge of some of the foremost and upcoming specialists in Finland outside my own restricted field. At the postgraduate seminars organized by the Department of Finnish language and literature, which I had the pleasure of attending twice, I was happy to receive from Pentti Leino both common-sense and worthwhile feedback from the Fennistic perspective on interpreting quantitative results arising from my statistically-oriented line of research. When I have needed guidance concerning the etymological roots of the THINK verbs that I study or help with rarer infinitival or other infrequent forms in Finnish, I have had the opportunity to consult Jussi Ylikoski's unfailing store of knowledge. I have also been able to indulge in my interest in writing systems by studying the cuneiform script, elementary Akkadian and Sumerian under the tutelage of Simo Parpola, and Maya hieroglyphs under the guidance of Harri Kettunen. These studies have provided me a truly welcome and worthwhile, culturally and historically rich counterpart and relaxation to my otherwise quite numeric study. Moreover, I wish to thank Harri for his peer-support during some trying times in my research process, which I hope to be able to reciprocate.

In the linguistic analysis of my research data, I was fortunate enough to have for several months two research assistants, Marjaana Välisalo and Paula Sirjola, to produce comparative analyses for parts of the data. I can only admire the thorough inquisition to which I was subjected concerning my principles of classification, which certainly became more uniform and explicit in the process, and I commend them both for the diligence in applying these principles in their own analysis of the data.

In the numerous statistical analyses undertaken throughout my research process, I am grateful to Daniel Aalto (Helsinki University of Technology), Pertti Palo (Helsinki University of Technology), Lauri Tarkkonen, Kimmo Vehkalahti, Seppo Nyrkkö, and Kari T. Vasko (CSC – IT Center for Science) for the invaluable advice that I have received countless times. Especially indebted I am to Daniel Aalto, who inspected Sections 3-4 describing the array of statistical methods applied in this study, although any possible faults and mistakes naturally remain mine alone. I also want to thank Seppo Nyrkkö for his assistance in getting *R* to work properly with the Finnish character set on my portable Macintosh, and Jarno Tuimala and Yrjö Leino (both of CSC – IT Center for Science) for providing me all the programming and other assistance necessary for the parallel computations that I undertook on CSC's brand-new *murska* supercluster system.

Likewise, I owe many thanks to the administrative and support staff at our Department, in the Faculty of Arts, and in the LANGNET graduate school for keeping the paperwork and infrastructure running smoothly, so that I and other researchers could mostly concentrate on academic matters. Especially, I want to mention Johanna Ratia, Anders Enkvist, Jaakko Leino, Ulla Vanhatalo, Hanna Westerlund, Tiina Keisanen, Miroslav Lehecka, Jussi Luukka, and Kristiina Norlamo for doing their utmost to facilitate my research work, not to forget Kirsti Kamppuri, Ardy Pitré, and others who no longer work in the academic world.

In the final drive to get this text ready for publication and public defense, I am glad that I could count on the efforts of Leila Virtanen, Jack Rueter, Graham Wilcock, in order to ensure that the quality of the language matched up to the content of the text. In particular, I am indebted to Jack Rueter and Leila Virtanen for their willingness to help me iron out the last stylistic wrinkles at very short notice.

My research has been financially supported for seven years within the GILTA project (funded by *TEKES*, the National Technology Agency of Finland within its USIX technology program: 40943/99) and the LANGNET graduate school (funded by the Finnish Ministry of Education), for which I am most grateful. In addition, I wish to thank the Chancellor of the University of Helsinki for several travel grants. Moreover, I am thankful to Juhani Reiman and Simo Vihjanen at Lingsoft for making generous use of my prior experience and linguistic knowledge at a remunerative rate to supplement my income. Furthermore, I also appreciate the many discussions concerning the philosophy of language and its study that I have had with Simo Vihjanen. I am also thankful to Pertti Jarla, Punishment Pictures and PIB for permission to use a strip from *Fingerpori*.

Finally, though a level of seclusion and social fasting seems to facilitate the conclusion of a doctoral dissertation, it would be hard to follow through with the entire post-graduate project solely within an academic bubble, with no connections to the rest of the world. In this respect, I am deeply grateful for the unwavering encouragement and patient understanding by my old and newer friends from Kulosaari, Helsinki University of Technology, TKY, Prodeko, the Polytech Choir, Teekkarispeksi, Kårspex, Spock's HUT, Aspekti, SKY – Linguistic Association of Finland, Flemarin Piikki, Pikku-Vallila, FC Heinähattu, Rytmi, and whatnot. Special mentions are awarded to the *Antti Club* for carrying on our culinary circuit, and *Ryhmä Rämä* 'Team Daredevil' for keeping in touch despite our distributed locations around the globe. I also thank my two godchildren, Juuso Rantanen and Otto Jalas, for the instant getaways I have always had in their company, providing an excellent antidote when I had lost myself too deep into scientific thought.

Last but not least, I wish to express my profound gratitude and appreciation to my immediate family for all the love, care, and concrete support that they have provided me throughout this longer-than-originally-anticipated and demanding process. To mention only a few details, I am thankful to my brother Tuomas Arppe for furnishing me a practical field stretcher to my workroom so that I need not take an afternoon nap on the cold floor, and to my sister Helena Arppe for helping me arrange the stylistic review of my dissertation text. What is more, I am immensely indebted to my father and mother, Juhani and Raija Arppe, for understanding and backing whatever decisions I have taken in my life, even supporting the very final phase financially, and in encouraging me to boldly go for my own thing.

Antti Arppe
Helsinki
December 2008

Table of Contents

1	Introduction.....	1
1.1	The empirical study of language and the role of corpora	1
1.2	Synonymy and semantic similarity	7
1.3	Theoretical premises and assumptions in this study	13
1.4	The goals of this study and its structure.....	13
2	Data and its linguistic analysis.....	17
2.1	Selection of the studied synonyms.....	17
2.1.1	Background in prior studies	17
2.1.2	Screening out the interrelationships of COGNITION lexemes by their dictionary definitions	20
2.1.3	Clustering the COGNITION lexemes by statistical means	23
2.1.4	Extracting frequencies for THINK lexemes	25
2.2	Selection of the contextual features and their application in the analysis.....	28
2.3	Present descriptions of the studied THINK synonyms.....	33
2.3.1	THINK lexemes in Pajunen’s <i>Argumenttirakenne</i>	33
2.3.2	THINK lexemes in <i>Suomen kielen perussanakirja</i> and <i>Nykysuomen sanakirja</i>	36
2.3.3	The etymological origins of the selected THINK lexemes.....	48
2.4	The compilation of the research corpus and its general description	50
2.4.1	General criteria.....	50
2.4.2	Main characteristics of the final selected research corpus.....	52
2.4.3	Coverage of contextual features in the research corpus.....	55
2.4.4	Representativeness, replicability, reliability, and validity	66
2.5	Preparation of the corpus for statistical analysis.....	69
3	Selection and implementation of statistical methods.....	71
3.1	Practical implementation of the employed statistical methods.....	71
3.2	Univariate methods	73
3.2.1	Homogeneity or heterogeneity of the distribution – independence or dependence?	75
3.2.2	Measures of association	86
3.2.3	Grouped univariate analysis for a set of related contextual features	94
3.3	Bivariate methods	105
3.3.1	General considerations.....	105
3.3.2	Pairwise associations of individual features	106
3.3.3	Pairwise comparisons of two sets of related features	108
3.4	Multivariate methods	113
3.4.1	Logistic regression analysis with nominal outcomes and variables	113
3.4.2	Selection of variables in multivariate logistic regression	116
3.4.3	Alternative heuristics of multinomial regression analysis	119
3.4.4	Evaluating the polytomous logistic regression models and their performance	126
3.4.5	A detailed example of a (binary) logistic regression model	139
3.4.6	Other possible or relevant multivariate methods	150

4	Univariate and bivariate results	153
4.1	Univariate analyses	153
4.1.1	General results	153
4.1.2	Characterizations of the studied THINK lexemes on the basis of the univariate results	160
4.1.3	Comparison of the univariate results with existing lexicographical descriptions	163
4.2	Bivariate correlations and comparisons	170
4.2.1	Pairwise correlations of singular features	170
4.2.2	Pairwise associations of grouped features	182
5	Multivariate analyses	187
5.1	Selection of variables	187
5.2	Comparisons of the descriptive and predictive capabilities of the different heuristics and models	198
5.2.1	Comparing the various heuristics with respect to the full model.....	198
5.2.2	The lexeme-wise breakdown of the prediction results	201
5.2.3	Comparing the performance of models with different levels of complexity.....	204
5.3	Relative lexeme-wise weights of feature categories and individual features.....	209
5.3.1	Overview of the lexeme-specific feature-wise odds.....	209
5.3.2	Lexeme-wise analysis of the estimated odds	212
5.3.3	Feature-wise analysis of the estimated odds.....	213
5.3.4	Comparison with the univariate results.....	216
5.3.5	Comparison with descriptions in current dictionaries	217
5.4	Assessing the robustness of the effects.....	220
5.4.1	Simple bootstrap and writer-cluster resampling	220
5.4.2	Assessing the effect of incorporating extralinguistic features.....	224
5.5	Probability estimates of the studied lexemes in the original data.....	228
5.5.1	Overall assessment of lexeme-wise probability estimates.....	228
5.5.2	Profiles of instance-wise distributions of the lexeme-wise probability estimates.....	237
5.5.3	“Wrong” choices in terms of lexeme-wise estimated probabilities.....	239
5.5.4	Contexts with lexeme-wise equal probability estimates – examples of synonymy?.....	241
5.5.5	Deriving general scenarios of probability distribution profiles with clustering.....	246
5.6	New descriptions of the studied synonyms.....	248
6	Discussion.....	253
6.1	Synonymy and its study and description in light of the results	253
6.2	Hypotheses for experimentation on the basis of the results.....	259
6.3	Suggestions for other further research and analyses.....	261
7	Conclusions.....	265
	Corpora	268
	References	269

Appendices⁴

Appendix A. Evaluation of the interchangeability of selected THINK lexemes among the example sentences provided in PS (Haarala et al. 1997)	288
Appendix B. Details concerning the selection of the studied THINK lexemes	289
Appendix C. Description of the various stages and levels in the linguistic annotation process and of the contextual features applied therein	317
Appendix D. List of morphological, surface-syntactic and functional syntactic features used in the linguistic analysis.....	350
Appendix E. Figures and selected details concerning the performance of the FI-FDG parser and the consistency of the manual annotation on the research data.....	384
Appendix F. Linguistic analyses of the lexical entries of the studied THINK lexemes in <i>Suomen kielen perussanakirja</i> (PS) and <i>Nykysuomen sanakirja</i> (NS).....	422
Appendix G. Lexeme-wise aggregates of the linguistic analyses of the lexical entries for the studied THINK lexemes, integrating the contents of both <i>Perussanakirja</i> and <i>Nykysuomen sanakirja</i>	436
Appendix H. Posited etymologies of selected THINK lexemes	442
Appendix I. An in-depth discussion of the text types incorporated in the sources selected into the research corpus, as well as a detailed description of the compilation and resultant characteristics and composition of this corpus	444
Appendix J. Frequency data concerning the selected sources for the research corpus, namely, Helsingin Sanomat (1995) and SFNET (2002-2003), plus an original data sample from Helsingin Sanomat	466
Appendix K. A Zipfian alternative for scrutinizing expected distributions of features among the studied THINK lexemes.....	478
Appendix L. An in-depth discussion of the conceptual foundations and associated parameters for the measures of association presented in Section 3.1.2.....	488
Appendix M. Interaction of medium with person/number features, semantic types of agents, and semantic and structural types of patients, studied in a dichotomous model pitting <i>ajatella</i> against the other three THINK lexemes	499
Appendix N. A full-depth presentation and discussion of selected univariate results	502
Appendix O. Some general results concerning the Zipfian scrutiny of the distributions of the studied features among the selected THINK lexemes.....	542
Appendix P. Complete univariate results	544
Appendix Q. Complete bivariate results.....	560
Appendix R. Complete multivariate variable sets and results	568
Appendix S. Brief descriptions of the main <i>R</i> functions applied throughout this dissertation	597

⁴ Included in the PDF version of this dissertation to be found at
URL: <http://ethesis.helsinki.fi/>

1 Introduction

1.1 The empirical study of language and the role of corpora

Human language is a multimodal phenomenon, involving physical, biological and physiological, psychological and cognitive, as well as social dimensions. Firstly, language is physical through the sound waves, gestures, written symbols, and electronic forms with which it is communicated by one language user to another, and which are the manifestations of language that we can externally perceive and observe easily. Secondly, language is biological and physiological with respect to the organs and senses that produce, receive, and process the physical manifestations of language, including the vocal tract, ears and hearing, eyes and sight, hands, and in some rarer cases also touch, and most importantly, the brain. Thirdly, language is psychological and cognitive in that its externally observable manifestations are linked with a psychological representation in the human cognitive system, yielding the Saussurean dichotomy between form and meaning. Fourthly, language is a social phenomenon: such meaning – and even the associated form – is constructed through and as a part of the collective activity and interpersonal communication of human beings; with no communicative or other socially shared functions language and its manifestations are meaningless.

Therefore, it is surprising that the study of language, *linguistics*, has, at least in the second half of the twentieth century seen a predilection for methodological monism (cf. Hacking 1996: 65-66). Firstly, the influential generative school (Chomsky 1965: 4, 65, 201) in its various incarnations has traditionally deemed language use, as manifested in, for example, corpora, as deficient and erroneous evidence about language as a system and the rules and regularities that it consists of, and thus an unreliable source of evidence in the study of language. The result of this reasoning has been the elevation of *introspection*, prototypically by the linguist himself with respect to his own *intuitions* about a linguistic phenomenon, as the primary type of linguistic evidence. The associated marginalization of corpora by generativism, however, has not been as uniformly categorical as is generally assumed (Karlsson, forthcoming 2008). Nevertheless, from the perspective of science studies, this in effect methodological exclusiveness by generativists – Noam Chomsky in particular – has been criticized as simplistic, as though it were implying that there was only one proper “style” of conducting modern science, that is, *hypothetical modeling*, when there are in fact many appropriate methods (Hacking 1996: 64-66, see also Crombie 1981: 284 for the general typology of scientific “styles”⁵); similar critical arguments have also been voiced within linguistics (Chafe 1992: 96; Wasow and Arnold 2005: 1484). Furthermore, within the generative school itself as in others, intuition and introspection, specifically as undertaken by the researchers themselves, have been demonstrated to be unreliable and inconsistent as a method (e.g., Schütze 1996; Bard et al. 1996; Sampson 2001, 2005; Gries 2002; Wasow and Arnold 2005; Featherston

⁵ The varieties of “styles” of science according to Crombie (1981: 284) are “(1) simple postulation established in the mathematical sciences, 2) the experimental exploration and measure of more complex observable relations, 3) the hypothetical construction of analogical models, 4) the ordering of variety by comparison and taxonomy, 5) the statistical analysis of the regularities of populations and the calculus of probabilities, and 6) the historical derivation of development. The first three of these methods concern essentially the science of individual regularities, and the second three the science of the regularities of populations ordered in space and time.” [numbering added by A.A.]

2007⁶). Some of these critics even appear prepared to go as far as discrediting any type of linguistic research method building upon intuition, such as elicitation or experimentation, a view most strongly vocalized by Sampson (2001: 129; 2005: 17; 2007a: 15; 2007b: 119).

In the late twentieth century, the natural use of language, collected and compiled as *corpora*, has predominantly been presented as *the* empirical solution to the inadequacy of introspection. This development has been particularly strengthened by the increasing availability of texts in electronic form, attributable firstly to the extremely successful diffusion of the personal computer, and, more recently, by the rapid dissemination of the Internet and the accelerating expansion of its content. However, proponents of this source of evidence have also been inclined to methodological preferentialism if not outright monism, as Chafe (1992: 96) so aptly puts it. At the least, it is fair to say that many language researchers who identify themselves as “corpus linguists” would elevate corpora as the most preferred or the most precise source of linguistic evidence (e.g., Leech et al. 1994; Sampson 2001, 2005; Gries 2002: 38); some would even go as far as to rank recordings of spoken language – representing the most natural and basic mode of linguistic behavior – first in a hierarchy of linguistic data (Sampson 2001: 7-10).

However, even corpus linguists are willing to admit that corpora cannot always provide a satisfactory or complete answer to all linguistically interesting or important research questions. First of all, it is difficult – if not impossible – to study rarer linguistic phenomena on the basis of corpora alone, as it is hard to distinguish such infrequent items from genuine errors, slips of the tongue, or effects of linguistic/cognitive disorders in production, or yet unestablished new forms or constructions resulting from linguistic change in the making (e.g., Sampson 2007a: 14). What is more, the inability of corpora, being fundamentally samples of language use, to produce direct *negative evidence* has also traditionally been presented as a limitation to their status as linguistic evidence. In other words, the absence of a given linguistic phenomenon in some corpus, while it may be indicative of the rarity of this phenomenon and thus a low expected probability of occurrence, cannot be taken as definitive, conclusive evidence that the phenomenon in question would *not* with certainty be a proper and acceptable linguistic structure (e.g., Atkins and Levin 1995: 87, 108; Hanks 1996: 78).⁷ Nevertheless, the issue of negative evidence is a problematic one for any empirical science. Finally, there is an ongoing discussion concerning the representativeness of corpora and how to improve this state of affairs,

⁶ It is peculiar to note that in Featherston’s (2007: 271) judgment the original criticisms concerning self-introspection as evidence by Schütze and others, presented already more than a decade ago, have not yet been fully accepted among a significant number of the generative linguists; apparently “old habits die hard.”

⁷ This has been countered with the argument that natural sciences, such as physics, typically presented as the model to be followed in linguistics and other human-oriented, social sciences in order to be considered “proper”, “hard” sciences, generally do not require negative evidence. For instance, the fact that we always see bricks falling down and never levitating upwards (by themselves) does allow us to conclude both that (a) gravity causes objects to be attracted towards each other, as well as that (b) gravity does *not* cause objects to be repelled away from each other, unless some other force is at play; though one could logically counter that we may have missed observing (b) to happen somewhere in the universe (adapted from a posting by Geoffrey Sampson 4.1.2008 to the Corpora mailing list). This view of the natural sciences as not using any negative evidence at all has been disputed by, e.g., John Goldsmith (posting 3.1.2008 to the same list).

that is, what exactly constitutes the entire population of different types of language usage events that corpora currently represent, and in the future could and should incorporate, and what the proportions of these types sampled into corpora as well as their overall sizes should be (e.g., Clear 1992; Biber 1993; Váradi 2001; Kilgariff and Grefenstette 2003: 336, 40-341; Leech 2007: 132). Any foreseeable developments, however, do not eliminate the fact that even an extremely large and comprehensive – and thus a very representative corpus or set of corpora – would still be representative of only one, even if central, aspect of linguistic behavior, namely, usage, with a bias towards production.

In fact, the range of different types of phenomena that can be considered part of or relevant to linguistic behavior, and which thus can also provide us with linguistic evidence, is quite diverse (e.g., Penke and Rosenbach 2004b: 485). Besides introspection or corpora made up of written or spoken language, we can also solicit judgements by (typically) native speakers concerning the grammaticality or acceptability of a linguistic form or structure; this is at its core in fact more commonplace as a linguistic activity than one might initially expect, since we encounter or engage in it all the time when we teach or learn a language, or when we attempt to guess the dialect of a native speaker or the underlying mother tongue of a second-language learner. Such elicitation has long been established as the central method in creating descriptions for previously (in scientific terms) unstudied languages, often in practice from scratch, within the domain of field linguistics. In addition, we can use linguistic errors and slips of the tongue by “normal” language users, or errors committed by language learners or people with linguistic/cognitive disorders of various sorts. Furthermore, we can study reaction times to visual or oral linguistic stimuli, the speed and progress of reading and the associated eye-movements, thus linking linguistics to the methodology of psychology. Or we can use neuroimaging to study how various parts of the brain are activated in conjunction with linguistic stimuli and tasks, linking linguistics to biology and cognitive science. On closer inspection, we see that these other sources of evidence are customarily found outside the confines of theoretical “core” linguistics “proper” as it has been largely conceived in the second half of the twentieth century, in its many “hyphenated” subfields.

Penke and Rosenbach (2004b: 485-491) attempt to give some structure to this methodological and evidential multitude by providing a tentative general classification along three dimensions. The first of these concerns *quantitative* vs. *qualitative* focus, that is whether we are interested in how many times a linguistic phenomenon occurs (in comparison to others), or in the dichotomy whether it occurs at least once or not at all. The second dimension contrasts *direct* vs. *indirect* evidence, distinguishing the (mainly brain-specific) psycho-cognitive-physiological processes that produce and interpret language from the external manifestations we can easily observe. Last, the third dimension distinguishes *spontaneous* vs. *elicited* production, that is the linguistic behavior produced naturally, independently of the researcher, or via interviewing, questionnaires or experiments in controlled settings. Introspection, if understood as self-elicitation, would fall into the category of qualitative, indirect and elicited evidence, as I consider, like Penke and Rosenbach (2004b: 492, citing Schütze 1996), linguistic judgements – or *competence* in generative terminology – as a form of *performance* rather than as a “direct” channel to our internal, overall psychological representation of language as a system. In turn, corpora as they are

traditionally conceived would be categorized in this typology as either quantitative or qualitative indirect spontaneous evidence, thus concerning only two of the altogether six possible slots of evidence types. Therefore, the unquestionably multimodal and multidimensional nature of language would appear to quite naturally lead to a pluralistic conception of how language can – and ought to – be studied and explained in order to attain a comprehensive understanding of language as a phenomenon, rejecting methodological exclusiveness and monism, and consequently also the primacy of one type of evidence over the others, be it introspection or corpora (cf. Chafe 1992: 96; see also Gries et al. 2005a: 666).

Until quite recently, linguistic research seems to have been characteristically restricted to one or another single type of data and associated research method as the only source of evidence. In fact, it appears that only within the last decade has the discipline started to explore and exploit the combination of multiple data sources and multiple methods as evidence. Moreover, it is no longer uncommon to see two or even more evidence types and methods used within one study. For instance, out of a collection of 26 studies in Kepser and Reis (2005a), half (13) employed two, or in a few cases, even more⁸ different empirical data types and methods. Other good examples of single studies incorporating multiple methods have been undertaken by Gries (2002) and Rosenbach (2003) concerning the English possessive alternation, Gries, Hampe and Schönefeld (2005a, 2005b) concerning English *as*-predicative structures, Gries (2003b) concerning the English dative alternation, Featherston (2005) concerning a range of English and German syntactic structures, and Arppe and Järvikivi (2007b) concerning a Finnish synonymous verb pair.

A multimethodological perspective can also be achieved by applying a method previously unapplied to a research question for which evidence has already been derived previously with another method, often but not always by other researcher(s). Examples of this latter kind of research set-up are by Bresnan (2007), who tests with experimentation the corpus-based results from Bresnan et al. (2007) concerning the English dative alternation, Vanhatalo (2003, included also in 2005), who contrasts corpus-based results regarding a Finnish synonymous adjective pair by Jantunen (2001, included also in 2004) with questionnaire survey data, and Kempen and Harbusch (2005), who compare corpus-based analyses against experimental results concerning word order in the midfield (*Mittelfeld*) of German subordinate clauses reported by Keller (2000). As Kepser and Reis (2005b) point out, each data type and method increases our linguistic knowledge, not only by confirming earlier results from other data types but also by adding new perspectives to the picture. Although the benefits of such triangulation are obvious, the various mixtures of data from different sources of evidence, with different origins and characteristics, must also be adequately reconciled. Especially in such multimethodological research conducted by different researchers with independently undertaken analyses and independently reported results, the challenge resides in keeping the selected explanatory factors and their interpretation plus their practical operationalization as consistent and as explicit as possible from one study to the next. A key characteristic of all of the aforementioned multimethodological studies or methodological comparisons is that

⁸ The number of data types and methods is in many cases difficult to specify exactly, as a study might incorporate, possibly quite comprehensively, results from research undertaken by others, or as a single experiment may be analyzed via two, clearly distinct perspectives (e.g., linguistically cued visual target choice and associated eye-movement), thus providing two types of data.

their results are essentially convergent, over a range of languages as well as linguistic phenomena. Nevertheless, due to their distinct premises and analytical perspectives, this convergence does not render the different types of evidence as mutually redundant.

In combining the most common types of linguistic evidence in an effective manner, Gries (2002: 27-28) has suggested a general research strategy based on the individual strengths and weaknesses of the different methods, which Arppe and Järvikivi (2007b: 151-152) have extended and specified with respect to two commonly used types of experimentation, resulting in the relationships between the different types of evidence exemplified in Table 1.1 below. As their example case, Arppe and Järvikivi (2007b) studied the differences in the (written) usage and experimental judgements concerning two Finnish synonymous verbs, *mieltiä* and *pohtia* ‘think, reflect, ponder’, with respect to the main semantic subtypes of their subject/AGENT (INDIVIDUAL vs. COLLECTIVE) as well as related morphological person features (FIRST vs. THIRD PERSON). First, one should begin by constructing research hypotheses on the basis of earlier research and one’s own professional linguistic introspection concerning the selected object of research. Next, the thus formulated hypotheses can be fine-tuned and specified as well as roughly validated for further examination using both qualitative and quantitative analysis of the pertinent corpus data. At this stage, one can already be fairly confident that (in relative terms) frequent phenomena are also highly acceptable, but as the population of the evidence combination slots in Table 1.1 based on the results by Arppe and Järvikivi (2007b) indicate, rareness does not necessarily correlate with lower acceptability.

Table 1.1. Relationships between different types of evidence, namely, between frequencies from corpora and (forced-choice) preference and acceptability judgements from experiments (Table 5 from Arppe and Järvikivi 2007b: 152).

Preferred	Dispreferred	Frequency/ Judgement	Unacceptable	Acceptable
<i>mieltiä</i> + FIRST PERSON SINGULAR+ INDIVIDUAL <i>pohtia</i> + COLLECTIVE (THIRD PERSON SINGULAR)	∅	Frequent	∅	<i>mieltiä</i> + FIRST PERSON SINGULAR+ INDIVIDUAL <i>pohtia</i> + COLLECTIVE (THIRD PERSON SINGULAR)
∅	<i>pohtia</i> + FIRST PERSON SINGULAR+ INDIVIDUAL <i>mieltiä</i> + COLLECTIVE (THIRD PERSON SINGULAR)	Infrequent	<i>mieltiä</i> + COLLECTIVE (THIRD PERSON SINGULAR)	<i>pohtia</i> + FIRST PERSON SINGULAR+ INDIVIDUAL

Consequently, one can use experimentation to get a better understanding of the (relatively) rarer phenomena. If one is interested primarily in usage preferences, for example, selections among alternatives, forced-choice experiments are the method of choice, as they would appear to roughly correlate with corpus-based results. In comparison, acceptability judgement ratings are a more precise and explanatorily

more powerful method of experimentation, as they are able to bring forth subtle yet significant differences among alternatives, even when none are evident in either corpora or through the forced-choice tasks (Featherston 2005). While corpus-derived results of alternative possible linguistic structures would tend to follow a Zipfian distribution (Zipf 1935, 1949), with the best-judged alternative also occurring with the relatively highest frequency, but the rest with very few if any occurrences at all, acceptability ratings of the same structures form a steadily declining linear continuum from the best-judged to the lowest-judged items, as is evident in Figure 1.1 below (Featherston 2005). Furthermore, according to Featherston, there would not appear to be any significant discontinuities among the range of alternative structures which would clearly divide them into grammatical and ungrammatical ones, the latter a view suggested by Kempen and Harbusch (2005) as well as Sorace and Keller (2005), with which I myself would be inclined to disagree.

Nevertheless, from the overall perspective, acceptability ratings do not contradict corpus-based frequencies or selections in forced-choice tasks, since the best-judged alternatives are also relatively more frequent in corpora than the worse-judged alternatives. However, this relationship between the two types of evidence is asymmetrical because relative rareness does not directly imply acceptability; a rare item can be either fully acceptable or clearly unacceptable; nor does acceptability directly imply a high frequency (Arppe and Järvikivi 2007b). Nonetheless, as this earlier study was restricted to only two alternative synonyms and to only a few – though important – contextual features, my intention in this dissertation is to extend the scope of study to encompass both the number of alternative lexemes and the range of contextual features considered (to be revisited in Section 1.2).

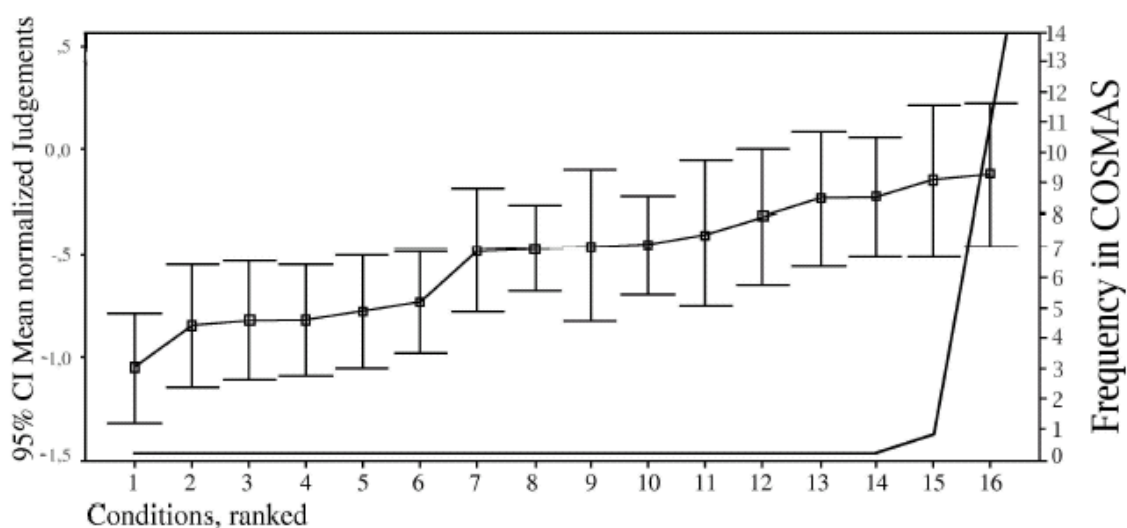


Figure 1.1. The contrast between corpus (COSMAS⁹) frequency data and experimental judgement data on the same phenomenon (corresponding to Figure 1 in Featherston 2004: 52, and Figure 4 in Featherston 2005: 195).

⁹ The acronym COSMAS stands for Corpus Search, Management and Search system, which gives online access to the German language corpora of the Institut für Deutsche Sprache (IDS) in Mannheim, Germany, exceeding currently well over one billion words; URL: <http://www.ids-mannheim.de/cosmas2/>.

Nevertheless, in comparison to the processing of corpora, experiments are considerably more time-consuming and laborious as well as subject to factors beyond the researcher's personal control – after all, they require a substantial number of committed informants in a specified setting, perhaps also with special measurement instruments in order to produce scientifically valid results. Thus, using corpus-based analysis first to prune and select only a small set of the most relevant or otherwise interesting hypotheses for further testing with focused experimentation is well motivated on practical and economic grounds. Furthermore, despite the deficiencies of introspection as a primary source of evidence, a researcher can, and in fact has to use his linguistic intuition and introspection to interpret the results and to adjust the research hypotheses throughout the different stages and associated methods in the aforementioned research cycle (Cruse 1986: 10-11; Sampson 2001: 136-139; cf. “heuristic exploratory device” in Gries 2002: 27).

As a final note, the concepts *evidence*, *data* and *method* are often used in an overlapping manner and may thus be difficult to clearly distinguish from one another. For instance, in the case of a corpus-based study, one could regard a corpus as the raw *data*, and *evidence* as simply various snippets selected from the corpus pertaining to the studied linguistic phenomenon, or in a varyingly more complex form the frequencies of various selected individual linguistic items and of their co-occurrences extracted from the corpus, and whatever analysis one can and might perform on this frequency information. As for what constitutes the *method*, one could, in the simplest case, consider making direct observations from a given corpus as the method; in the more complex analyses the observations would be based on a sequence of a variety of non-trivial tasks starting with the collection or selection of a corpus, its linguistic annotation, and the choice of appropriate statistical methods, and so on. So, in a sense, a corpus can play the role of raw data, of method and of evidence.

1.2 Synonymy and semantic similarity

The linguistic phenomenon studied in this dissertation is lexical *synonymy*, which I understand as *semantic similarity* of the nearest kind, as discussed by Miller and Charles (1991), that is, the closest end on the continuum of *semantic distance* between words. My general theoretical outlook is therefore linguistic empiricism in the tradition of Firth (1957), with meaning construed as *contextual*, in contrast to, for example, formal (de)compositionality (see, e.g., Cruse 1986: 22, Note 17; or Fellbaum 1998b: 92-94 for an overview of the relevant theories; or Zgusta 1971: 27-47 for the classical lexicographical model employing the concepts *designation/denotation*, *connotation*, and *range of application*). Thus, I operationalize synonymy as the highest degree of mutual substitutability (i.e., interchangeability), without an essential change in the perceived meaning of the utterance, in as many as possible in a set of relevant contexts (Miller and Charles 1991; Miller 1998). Consequently, I do not see synonymy as dichotomous in nature, but rather as a continuous characteristic; nor do I see the associated comparison of meanings to concern truth values of logical propositions or conceptual schemata consisting of attribute sets. In these respects, the concept (and questionability of the existence) of *absolute synonymy*, that is, complete interchangeability in all possible contexts, is not a relevant issue here

Nevertheless, it is fair to say that I regard as synonymy what in some traditional approaches, with a logical foundation of meaning, has rather been called *near-synonymy* (or *plesionymy*), which may contextually be characterized as “synonymy relative to a context” (Miller 1998: 24). However, like Edmonds and Hirst (2002: 107, Note 2), I see little point in expatiating on how to distinguish and differentiate synonymy (in general), near-synonymy, and absolute synonymy, especially since the last kind is considered very rare, if it exists at all. This general viewpoint is one which Edmonds and Hirst (2002: 117) ascribe to lexicographers, with whom I am inclined to align myself. A recent approach to synonymy which in my mind conceptually fleshes out the essence of this lexical similarity can be found in Cruse (2000: 156-160, see also 1986: 265-290), where synonymy is “based on empirical, contextual¹⁰ evidence”, and “synonyms are words 1) whose semantic similarities are more salient than their differences, 2) that do not primarily contrast with each other; and 3) whose permissible differences must in general be either minor, backgrounded, or both”.

In the modeling of the lexical choice among semantically similar words, specifically near-synonyms, it has been suggested in computational theory that (at least) three levels of representation would be necessary to account for fine-grained meaning differences and the associated usage preferences, namely, 1) a conceptual-semantic level, 2) a subconceptual/stylistic-semantic level, and 3) a syntactic-semantic level, each corresponding to increasingly more detailed representations, that is, granularity, of (word) meaning (Edmonds and Hirst 2002: 117-124). In such a model of language production (i.e., generation), synonyms are grouped together as initially undifferentiated clusters, each associated with individual coarse-grained concepts at the topmost level (1), according to a (possibly logical) general ontology. The individual synonyms within each cluster all share the essential, core denotation of the associated concept, but they are differentiated *in contrast to* and *in relation to each other* at the intermediate subconceptual level (2), according to peripheral denotational, expressive and stylistic distinctions, which can in the extreme be cluster-specific and fuzzy, and thus difficult if not impossible to represent simply in terms of absolute general features or truth conditions. Consequently, a cluster of near-synonyms is nonetheless internally structured in a meaningful way, which can be explicable, even if in a complex or peculiarly unique manner. By way of example, the expressive distinction can convey a speaker’s favorable, neutral or pejorative attitude to some entity involved in a discourse situation, while the stylistic distinction may indicate generally intended tones of communication such as formality, force, concreteness, floridity, and familiarity.

The last, syntactic-semantic level (3) in such a *clustered model of lexical knowledge* concerns the combinatorial preferences of individual words in forming written sentences and spoken utterances, for example, syntactic frames and collocational relationships. Though Edmonds and Hirst (2002: 139) do recognize that this level is in a complex interaction with the other two, they leave this relationship and the

¹⁰ One should note, however, that Cruse’s (1986: 8-10, 15-20) conception of *contextual relations* as the foundation of word meaning, and thus also synonymy, refers in terms of evidence rather to (the full set of) intuition-based judgments (possibly derived via experimentation) of the normality as well as the abnormality of a word in the totality of grammatically appropriate contexts, that is, including patterns of both *disaffinity* as well as *affinity*, and comparisons thereof, than the corpus-linguistic context of a word in samples of actually observed, natural language use (*productive output* in Cruse’s [1986: 8] terms).

specific internal workings of this level quite open. Working within this same general computational model, Inkpen and Hirst (2006) develop it further by also incorporating the syntactic-semantic level in the form of simple collocational preferences and dispreferences, though their notion of collocation is explicitly entirely based on statistical co-occurrence without any of the more analytical linguistic relationships (Inkpen and Hirst 2006: 12); they foresee that such contextual lexical associations could be linked with the subconceptual nuances which differentiate the synonyms within a cluster (Inkpen and Hirst 2006: 35). This fits neatly with the view presented by Atkins and Levin (1995: 96), representatives of more conventional linguistics and lexicography, that even slight differences in the conceptualization of the same real-world event or phenomenon, matched by different near-synonyms, are also reflected in their syntactic (i.e., contextual) behavior.

In general, this aforementioned computational model also resembles psycholinguistically grounded models concerning the organization of the lexicon such as WordNet (Miller et al. 1990) to the extent that lexemes are primarily clustered as undifferentiated synonym sets (i.e., *synsets*) that are associated with distinct concepts (i.e., meanings), while semantic relationships are essentially conceived to apply between concepts, signified in practice by the synsets as a whole. However, the WordNet model fundamentally considers all lexemes belonging to such individual synsets as mutually semantically equivalent, effectively ignoring any synset-internal distinctions that might exist among them (Miller et al. 1990: 236, 239, 241; Miller 1995; Miller 1998: 23-24, Fellbaum 1998a: 9).

Returning to the syntactic-semantic level, it has been shown in (mainly) lexicographically motivated corpus-based studies of actual lexical usage that semantically similar words differ significantly as to 1) the lexical context (e.g., English adjectives *powerful* vs. *strong* in Church et al. 1991), 2) the syntactic argument patterns (e.g., English verbs *begin* vs. *start* in Biber et al. 1998: 95-100), and 3) the semantic classification of some particular argument (e.g., the subjects/agents of English *shake/quake* verbs in Atkins and Levin 1995), as well as the rather style-associated 4) text types or registers (e.g., English adjectives *big* vs. *large* vs. *great* in Biber et al. 1998: 43-54), in which they are used. In addition to these studies that have focused on English, with its minimal morphology, it has also been shown for languages with extensive morphology, such as Finnish, that similar differentiation is evident as to 5) the inflectional forms and the associated morphosyntactic features in which synonyms are used (e.g., the Finnish adjectives *tärkeä* and *keskeinen* ‘important, central’ in Jantunen 2001, 2004; and the Finnish verbs *mieltiä* and *pohtia* ‘think, ponder, reflect, consider’ in Arppe 2002, Arppe and Järvikivi 2007b; see also an introductory discussion concerning inflectional distinctions of synonyms in general in Swedish, Danish, and Norwegian Bokmål in Arppe et al. 2000).

Recently, in their studies of Russian near-synonymous verbs denoting TRY as well as INTEND, Divjak (2006) and Divjak and Gries (2006) have shown that there is often more than one type of these factors simultaneously at play, and that it is therefore worthwhile to observe all categories together and in unison rather than separately one by one. Divjak and Gries (2006, forthcoming) dub such a comprehensive inventory of contextual features of a word as its *Behavioral Profile*, extending this notion to cover not only complementation patterns and syntactic roles as proposed by Hanks (1996),

who originally coined the concept, but *any* linguistic elements, whether phonological, morphological, syntactic, semantic, or other level of linguistic analysis, which can be observed within the immediate sentential context, adapting here the notion of the so-called *ID tags* presented by Atkins (1987).¹¹ Furthermore, Divjak and Gries also present one possible way of operationalizing and compactly quantifying this concept for each word as one co-occurrence vector of within-feature relative frequencies. In my mind, one could alternatively refer to this concept as the *Contextual Profile* or *Distributional Profile* of a word, as its primary components are the occurrences and distributions of linguistically relevant items or characteristics (or their combinations) which can be explicitly observed in a word's context in (a sample of) language usage. As noted earlier above, though Cruse's (1986: 8-10, 15-20) concept of *contextual relations* is quite similar in both name and intended purpose in defining linguistic meaning, it fails to examine explicitly the individual elements in the context itself.

All of these studies of synonymy have focused on which contextual factors differentiate words denoting a similar semantic content. In other words, which directly observable factors determine which word in a group of synonyms is selected in a particular context. This general development represents a shift away from more traditional armchair introspections about the connotations of and ranges of application for synonyms (e.g., Zgusta 1971), and it has been made possible by the accelerating development in the last decade or so of both corpus-linguistic resources, that is, corpora and tools to work them, such as linguistic parsers, and statistical software packages.

Similar corpus-based work has also been conducted on the syntactic level concerning *constructional alternations* (referred alternatively to as *synonymous structural variants* in Biber et al. 1998: 76-83), often from starting points which would be considered to be anchored more within general linguistic theory. Constructional alternations do resemble lexical synonymy, for the essential associated meaning is understood to remain largely constant regardless of which of the alternative constructions is selected; however, they may differ with respect to a pragmatic aspect such as focus. Relevant studies concerning these phenomena have been conducted by Gries (2002) and Rosenbach (2003) with respect to the English possessive constructions (i.e., [*NP*_{POSSESSED} of *NP*_{POSSESSOR}] vs. [*NP*'_{POSSESSOR} *NP*_{POSSESSED}]), Gries (2003a) concerning the English verb-particle placement, (i.e., [*VP* *NP*_{DIRECT_OBJECT}] vs. [*V* *NP*_{DIRECT_OBJECT} *P*]), and Gries (2003b) as well as Bresnan et al. (2007) concerning the English dative alternation, (i.e., [*GIVE* *NP*_{DIRECT_OBJECT} *PP*_{INDIRECT_OBJECT}] vs. [*GIVE* *NP*_{INDIRECT_OBJECT} *NP*_{DIRECT_OBJECT}]). The explanatory variables in these studies have been wide and varied, including phonological characteristics, morphological features and semantic classifications of relevant arguments, as well as discourse and information structure. With regard to Finnish, a good example of a syntactic alternation are the two comparative constructions, (i.e., [*NP*_{PARTITIVE} *A*_{COMPARATIVE}] vs. [*A*_{COMPARATIVE} *kuin* *NP*]), for example, *Pekkaa parempi* vs. *parempi kuin Pekka* 'better than Pekka', which is described by Hakulinen et al. (2004: 628-630 [§636-§637]), and prescriptively scrutinized by Pulkkinen (1992 and references). These two alternative constructions last mentioned are cross-linguistically well known and studied, and are

¹¹ Such an omnivorous attitude with respect to analysis levels and feature categories is an integral characteristic in machine learning approaches within the domain of computational linguistics.

considered to represent two distinct types in language-typological classifications (e.g., Stassen 1985, 2005).

With the exception of Gries (2002, 2003a, 2003b), Rosenbach (2003), Bresnan et al. (2006), Divjak (2006), and Divjak and Gries (2006), the aforementioned studies have in practice been monocausal, focusing on only one linguistic category or even a singular feature within a category at a time. Though Jantunen (2001, 2004) does set out to cover a broad range of feature categories and notes that a linguistic trait may be evident at several different levels of context at the same time (2004: 150-151), he does not quantitatively evaluate their interactions. Bresnan et al. (2006) have suggested that such reductive theories would result from pervasive correlations in the available data. Indeed, Gries (2003a: 32-36) has criticized this traditional proclivity for monocausal explanations and has demonstrated convincingly that such individual univariate analyses are insufficient and even mutually contradictory. As a necessary remedy in order to attain scientific validity in explaining the observed linguistic phenomena, he has argued forcefully for a holistic approach using multifactorial setups covering a representative range of linguistic categories, leading to the exploitation of multivariate statistical methods. In such an approach, linguistic choices, whether synonyms or alternative constructions, are understood to be determined by a *plurality* of factors, in *interaction* with each other. More generally, this can in my mind be considered a *non-modular* approach to linguistic analysis. Nevertheless, the applicable multivariate methods need to build upon initial univariate and bivariate analyses.

Furthermore, as has been pointed out by Divjak and Gries (2006), the majority of the above and other synonym studies appear to focus on word pairs, perhaps due to the methodological simplicity of such a setup; the same criticism of limited scope also applies to studies of constructional alternation, including Gries' own study on English particle placement (2003a). However, it is clearly evident in lexicographical descriptions such as dictionaries that there are often more than just two members to a synonym group, and this is supported by experimental evidence (Divjak and Gries, forthcoming). Though full interchangeability within a synonym set may *prima facie* be rarer, one can very well assume the existence of contexts and circumstances in which any one of the lexemes could be mutually substituted without an essential change to the conveyed meaning. Consequently, the differences observed between some synonymous word pair might change or relatively diminish when studied overall in relation to the entire relevant synonym group. This clearly motivates a shift of focus in synonym studies from word pairs to sets of similar lexemes with more than two members, an argument which has already been expressed by Atkins and Levin (1995: 86).

Finally, Bresnan (2007, see also 2006) has suggested that the selections of alternatives in a context, that is, lexical or structural outcomes for some combinations of variables, are generally speaking probabilistic, even though the individual choices in isolation are discrete (see also Bod et al. 2003). In other words, the workings of a linguistic system, represented by the range of variables according to a theory, and its resultant usage would not in practice be categorical, following from exception-less rules, but rather exhibit degrees of potential variation which becomes evident over longer

stretches of linguistic usage.¹² These are manifested in the observed proportions of occurrence for one particular dichotomy of alternating structures, given a set of contextual features. It is these proportions which Bresnan (2007) et al. (2007) try to model and represent with logistic regression as estimated expected probabilities, producing the continuum of variation between the practically categorical extremes evident in Figure 1.2. Both Gries (2003b) and Bresnan (2007, et al. 2007) have shown that there is evidence for such probabilistic character both in natural language use in corpora as well as language judgements in experiments, and that these two sources of evidence are convergent. However, these studies, too, have concerned only dichotomous outcome alternatives.

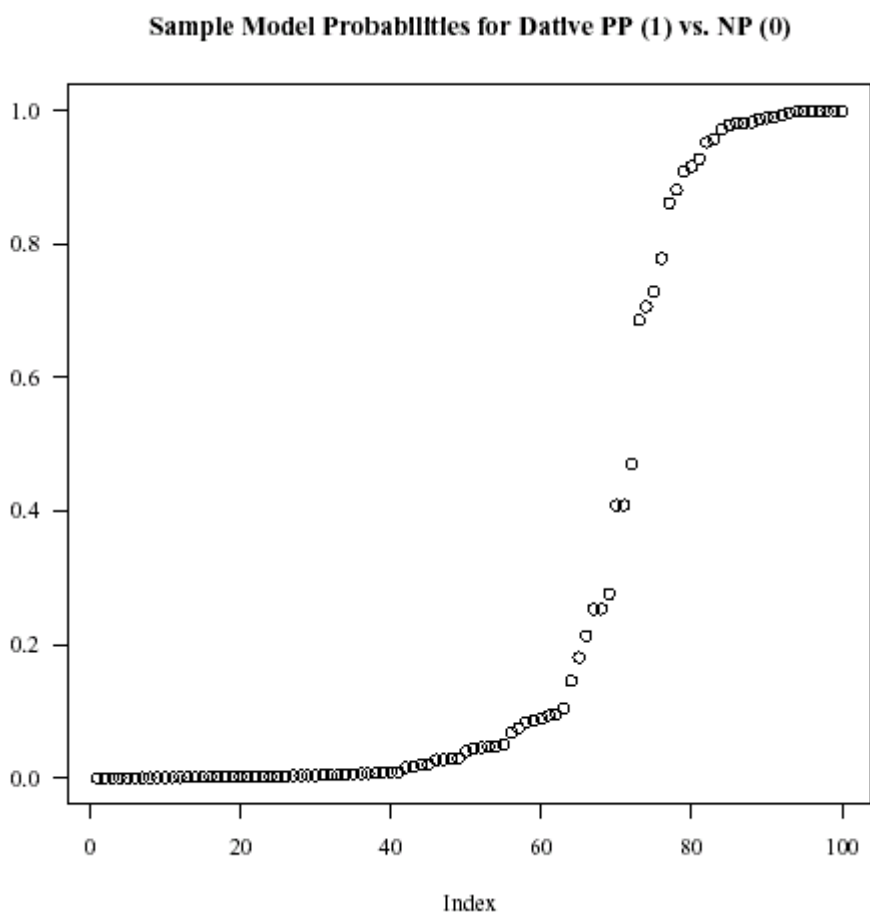


Figure 1.2. Sample of estimated expected probabilities for the English dative alternation (reproduced from Bresnan 2006: 4, Figure 1, based on results from Bresnan et al. 2007).

¹² From the perspective of understanding and explaining an empirical phenomenon, this means a shift from seeing *causes* as *deterministic*, “producing” an outcome by the “action of some universal and unfailing laws” to rather viewing *causes* as *probabilistic*, which merely “increase the likelihood” of such an outcome (Hacking’s [1981: 113] interpretation of Fagot’s [1981] discussion concerning the causes of death within medicine).

1.3 Theoretical premises and assumptions in this study

The key theoretical concepts characterizing this dissertation are 1) *contextuality*, 2) *synonymy*, and 3) *non-modularity* (or *constructionality*). Firstly, I conceive meaning as fundamentally contextual in the spirit of Firth (1957), rather than as compositional or propositional. Thus, I assume a strong interconnection between (contextual) distributional and semantic similarity (cf. Divjak and Gries 2006: 27-28). However, as I am working with Finnish, a language having a relatively flexible word order, I operationalize context as grammatical rather than linear, that is, within a fixed window of words. Secondly, I take a naïve, lexicographical view with regard to synonymy (cf. Cruse 1986: 265; Edmonds and Hirst 2002: 116). That is, I take it as granted that the lexicon contains pairs or sets of words which are mutually interchangeable without an essential change in their core conveyed meaning. Nevertheless, I suspect that this mutual interchangeability is not, and need not, be categorical. In terms of the contextual view, synonymy consequently means for me interchangeability in many/most contexts, but not necessarily all. Thirdly, I am a non-modularist with respect to linguistic analysis. In my view, regularities in co-occurrence and structure can be observed along an ungraded continuum with increasingly more abstract levels, from fixed collocations concerning words and morphemes up to conventional syntactic structures composed of general parts of speech or semantic classes of words. Nevertheless, we will often find it practical to segment linguistic analysis into several conventional levels such as morphology, syntax and semantics, but this does not to my mind require or lead to treating these levels as fully independent of, *autonomous* of, and detached from each other. This attitude towards linguistic analysis can be viewed as *constructional*, following Construction Grammar (e.g., Croft 2001); however, I will not adhere exclusively to any single linguistic theory in this dissertation.

1.4 The goals of this study and its structure

This dissertation is first and foremost a methodological study. My specific goal is to extend the study of one type of linguistic alternation, that is, synonyms and their choice, from dichotomous to polytomous settings, involving the lexical choice among more than two alternatives. Likewise, my purpose is also to move from simplistic, univariate explanatory models of lexical choice to more complex but powerful multivariate ones, explaining the phenomenon in question with a broad, representative range of linguistic and extralinguistic variables. Therefore, as a follow-up to Arppe and Järvikivi (2007b), the other members of the Finnish THINK synonym group with frequencies within the same (relatively high) magnitude as the original pair, have been selected for scrutiny in this study, resulting in the lexeme set *ajatella*, *mieltiä*, *pohtia* and *harkita*. Furthermore, instead of considering only the morphological structure or one specific argument of the studied verbs, as was the case in the former study, the range of contextual features included in the analysis in this dissertation will be extended to cover the entire syntactic argument structure of the studied verbs and the semantic subclassifications of the individual argument types, including extralinguistic features such as text type and register.

In terms of linguistic description, my purpose is to demonstrate how lexemes in a synonym group are used differently on the basis of their context, and to lay out the

individual features which determine these differences as well as assess their relative importance and weights. Ultimately, this should result in better descriptions of the lexemes I have chosen to study, and by way of replication, the overall improvement of broad-coverage lexicographical descriptions such as dictionaries. As for the development of linguistic research methodology, in this dissertation I intend to explain in detail why I have selected particular statistical methods and how these are in practice applied to corpus data, and, using examples, demonstrate how the results can be interpreted, thus also showing what the selected statistical methods have to offer from the linguistic perspective. This will be an exploratory study, with no specific hypotheses to prove or disprove other than the general assumption of the existence of differences in contextual preferences and dispreferences among the studied THINK lexemes, whatever the actual features involved may turn out to be. Consequently, the number of different individual features and feature categories covered in this study will be substantial.

On the general level, this setup of including a wide range of different features as well as a broader set of synonyms than a simple pair largely resembles that of Divjak and Gries (2006). However, my focus is rather on discovering features distinguishing the members of a synonym group from each other than on the internal grouping of a synonym set that these features also reveal. Furthermore, I will take the groupings of synonyms mostly as they are given in current, authoritative Finnish dictionaries, and I will not delve deeper into what corpus-based contextual evidence could indicate in this respect. Thus, of the three main challenges that Divjak and Gries (2006: 24-27) present for the study of synonyms, namely, 1) the lexicon-general *delineation* of words into synonym groups, 2) the internal *structuring* of these groups, and 3) the *description* of features which differentiate the individual synonyms from each other, my dissertation will focus on the third and last task. As Divjak and Gries (2006: 23-24) note, synonymy has received within Western linguistics far less attention in linguistic research than other lexical relationships such as polysemy. Though this “neglect” has to some extent been remedied by a string of recent studies both in Finland (e.g., Jantunen 2001, 2004; Arppe 2002; Arppe and Järviö 2007b; Vanhatalo 2003, 2005; Päiviö 2007) and abroad (e.g., Edmonds and Hirst 2002; Inkpen and Hirst 2006; Divjak and Gries 2006), the topic is neither far from exhausted nor conclusively resolved. Thus, one general objective of this study is to continue this trend as a worthwhile linguistic research topic and contribute to our understanding of synonymy as a linguistic phenomenon.

The long-term purpose of my research is to increase our understanding of what the relationship is between naturally produced language, that is, corpora, and the posited underlying language system that governs such usage. This concerns both 1) the use and choice among lexical and structural alternatives in language and 2) the underlying explanatory factors, following some theory representing language as a comprehensive system. A subsequent subservient methodological objective is how this can be modeled using various statistical methods with different levels of complexity, ranging from the univariate to the multivariate levels. A question which may at first glance appear secondary but which will turn out to be of general theoretical import is to what extent we can describe the observed variation in terms of the selected analytical features that conventional linguistic theory incorporates and works upon.

The structure of this study is as follows. In Section 2, I will begin by presenting the various linguistic aspects of this study, including the selection of the studied synonyms, the general principles of the various levels of linguistic analysis which are to be applied to the research corpus, and a description of the compilation as well as the characteristics of the selected research corpus. The individual levels of analysis and the associated specific features are fully presented in Appendix C. Furthermore, I will tease out what current lexicographical descriptions of the selected synonyms reveal about their usage, whether explicitly or (mostly) implicitly, in terms of the same features of linguistic analysis applied in this study, in order to provide a benchmark against which to compare the corpus-based results to follow later on. In Section 3, I will move on to lay down the battery of statistical methods to be used in this dissertation, starting off with several univariate methods, extending some of these to bivariate analysis, and finishing with multivariate methods, concentrating on polytomous logistic regression. Among the univariate methods, I will first address the assessment of the overall homogeneity/heterogeneity of a distribution and various follow-up tests, and then a more comprehensive exposition of various measures of association for nominal (non-ordered categorical) variables. The bivariate analysis methods will in fact be applications of the very same measures of association presented among the univariate methods. In addition, I discuss extending the method presented by Gries (2003a) from dichotomous to polytomous settings. Regarding the multivariate analysis methods, I will conclude by also presenting various ways of assessing the robustness and generalizability of the forthcoming results, with a major focus on bootstrapping procedures of several kinds. Throughout this Section, I will use only a limited set of linguistic analysis features, typically anchored in previous research if such exists, to take a detailed walk through the successive stages of the actual calculations, not to mention the various ways in which the thus retrieved results can be abstracted and interpreted from a linguistic perspective.

In Section 4, I will begin to show the results of applying these various statistical methods using the selected linguistic analysis features, again proceeding from univariate through bivariate to multivariate results. Here, I will present only the very last, abstracted end of the underlying statistical calculations, which will often be the result of several layers of simplification and interpretation. For those interested in specific analyses by each level and feature categories or the actual underlying figures, these can be found in the Appendices, while the details of their calculation have been exemplified earlier in Section 3. Among the univariate results presented in Subsection 4.1, I will begin with general observations of the feature-specific values of the selected measures of association, concluding with an attempt at *post hoc* generalizations presented in Subsection 4.1.2 and, finally, a comparison of the univariate results with existing lexicographical descriptions in Subsection 4.1.3.

The bivariate results discussed in Subsection 4.2 will most importantly prepare the ground for the considerable pruning down of the feature variables selected for the multivariate analyses covered in Section 5. Here, I will compare the performance of different techniques for implementing polytomous logistic regression for the full selected variable set as well as the explanatory power of different subsets of the selected variables representing varying degrees of analytical depth or intricacy. This will be followed by a full exposition of the feature-wise weights, that is, the estimated odds, for the finally selected feature set, and then an assessment of the robustness of the results with several techniques. Next, I will discuss at length the probability

estimates that polytomous logistic regression conveniently produces for each lexeme with respect to any combination of features incorporated in the developed model. This general Section 5 discussing the results will end in Subsection 5.6 with a suggestion for the new description of the studied synonyms building specifically upon the multivariate results, to be compared with the current ones analyzed earlier in Subsection 2.3.

For the most part, I will in Sections 4 and 5 link and discuss the specific observations and conclusions reached in this study with regard to previous research at those points when they are first encountered and presented, leaving only the most general conclusions to Section 6. In this next-to-last section, I will also sketch hypotheses for later experimentation, not to mention other ensuing areas for further research, many of which I have had to exclude from this dissertation due to limitations of space. Finally, a short overall conclusion highlighting the main results of this dissertation will be presented in Section 7. For a linguistically minded reader, either Section 4.1.2 presenting the *post hoc* general characterizations of the studied THINK lexemes, or Section 5.6 laying out the new lexicographical description scheme resulting from the multivariate analysis results, with *pohtia* as an example, might be the best starting points for getting a concise overview of what this dissertation is about in lexicographical terms.

In addition to the main text, the Appendices contain a wealth of information that I believe will be of genuine value and interest, but which is not integral to the central objectives of this dissertation. Moreover, a comprehensive collection of the data used in this study, the research corpora and their linguistic analyses, the scripts with which they have been processed and the tailored functions with which the statistical analyses have been undertaken in the *R* statistical computing environment, as well as the complete results, can all be found in the *amph* microcorpus under the auspices of CSC Scientific Computing,

<URL: <http://www.csc.fi/english/research/software/amph>>.

2 Data and its linguistic analysis

2.1 Selection of the studied synonyms

2.1.1 Background in prior studies

The set of four synonymous THINK lexemes scrutinized in this study, *ajatella*, *mieltiä*, *pohtia*, and *harkita*, were first and foremost selected because I had in earlier studies on my own and in co-operation with others focused extensively on one pairing among them, namely, *mieltiä* vs. *pohtia*, which I had considered semantically the closest ones of the group (Arppe 2002; Arppe and Järvikivi 2002; Arppe and Järvikivi 2007b). Although pairwise comparisons of synonyms are by far the most common in linguistics, perhaps in part because it is methodologically the easiest setting to pursue, synonymy as a phenomenon is by no means restricted to word pairs either conceptually or in practice, nor should its study be limited to such pairs, as Divjak and Gries (2006) argue and demonstrate, nor. For instance, we can find in dictionaries and thesauri often more than one synonym provided for many of the lexical entries. From my own experience, however, I must concede that in the case of most synonym sets, one can without much difficulty come up with contexts or connotations that clearly distinguish individual synonyms from the rest, often leaving one with only one pair (or pairs) which at least superficially are not immediately distinguishable from each other on the basis of one's professional (and native speaker's) linguistic intuition.

The original selection of the *mieltiä-pohtia* synonym pair and the entire synonym set of THINK verbs to which they belong was based on a rigorous process with the purpose of identifying lexemes for which their syntactic and semantic valency profiles as well as the “contamination” effect from their possible polysemous senses, and even extralinguistic factors such as their relative frequencies, should be as similar as possible (due to the *frequency effect*¹³ in linguistic processing, for which an overview is presented in Ellis 2002). The ultimate goal was thus to ensure *a priori* a degree of interchangeability as high as possible in the observable contexts, as a proxy for the nearest possible synonymy. Of course, one could have used a quantitative method such as the *sub-test* presented by Church et al. (1994) or its modification using exact statistics as suggested by Gries (2003c) to assess such factors empirically. However, because the present study, as well as my earlier ones, specifically use a corpus to uncover usage-based similarities and differences, I regarded other sources, independent of the chosen type of direct empirical evidence, as more appropriate.

These sources are the *Suomen kielen perussanakirja* in its various editions and forms (Haarala et al. 1994-1997, Haarala et al. 1997), that is, ‘The Standard Dictionary of Finnish’ hereafter denoted by the acronym *PS*, and the comprehensive descriptions of Finnish verbs by Pajunen (1982, 1988, 2001), which are all corpus-based, though each uses a different, yet in some cases overlapping, selection of source material. In addition, the predecessor to *PS*, *Nyky-suomen sanakirja* (Sadeniemi et al. [1951-1961] 1976), ‘Dictionary of Modern Finnish’ hereinafter denoted by the acronym *NS*, is also

¹³ Simply put, the frequency effect in the case of lexis means that the “recognition and processing of words is a function of their frequency of occurrence in the language” (Ellis 2002: 152); however, the underlying factors behind this empirical observation have been shown to be more complex than the simple definition would lead to believe, see, for example, Balota and Chumbley (1985), Schreuder and Baayen (1997), and Dahan et al. (2001).

consulted specifically in Section 2.3.2, which presents the extent to which the usage of studied lexemes has been described until now. The NS is a very comprehensive and extensive lexicographical work, which, exceeding 200,000 lexical entries, is almost twice the size of PS. However, it has essentially not been updated since it was compiled in 1929-1961 and is thus based on Finnish in the form it was used (and conceived to be) in the first half of the twentieth century. For this reason, I have primarily relied on the more up-to-date PS, as its contents have been during its existence since 1994 and thereafter under an on-going revision process by *Kotimaisten kielten tutkimuskeskus (KOTUS)*, the ‘Research Institute of the Domestic Languages in Finland’ <URL: <http://www.domlang.fi/>>, and even more so as PS in fact incorporates much of NS’s central content.¹⁴

In order to rule out pairs or sets of lexemes with potentially marked members resulting from relative infrequency, synonym candidates in the preceding studies were first ranked both by pairs and by entire sets according to the geometric averages of their relative frequencies (in the case of synonym sets considering only the values of their non-null members), based on a million-word corpus sample of Finnish newspaper text (a portion from Keski-suomalainen 1994),¹⁵ so that pairs or sets with high but at the same time also relatively similar frequencies came first. The synonym sets were extracted from the lexical database underlying the FINTHES inflecting thesaurus software module developed at Lingsoft.¹⁶ Using my own linguistic intuition as a native speaker of Finnish, I then scrutinized this ranking list from the top down in order to pick out promising candidates. In turn, these were evaluated in depth with respect to the similarity of their semantic and syntactic valency structures using both the specific descriptions by Pajunen (1982: 169, 180-182), when existent, and the definitions and usage examples from the lexical entries in PS (Haarala et al. 1997).

Regarding the first reference work, in its earlier form it covered explicitly only the most frequent or representative lexeme (or two) for each semantic field corresponding to a synonym group, in comparison to the more comprehensive coverage in its later, substantially revised extension (Pajunen 2001). This current version, however, had not yet appeared at the time of the initial research and it is still more exemplifying

¹⁴ Not to make things any less complicated, the latest versions of PS have in fact been marketed and distributed since 2004 under the name of *Kielitoimiston sanakirja* ‘The Dictionary of the Finnish Language Office’ (at KOTUS), denoted by the acronym *KS*. In terms of both its content and structure, *KS* is essentially an updated version of PS.

¹⁵ N.B. This subsample was part of the corpus used in my earlier studies (Arppe 2002; Arppe and Järvi-kivi 2007b), but not the part of the corpus material used in this study.

¹⁶ This Finnish synonym database was originally compiled by Katri Olkinuora and Mari Siirainen for Lingsoft at the behest of professor Kimmo Koskeniemi between 1989-1991 (see Arppe 2005b). This database has 7439 entries overall and approximately 29854 words (when not distinguishing multi-word synonyms as distinct units). These figures are certainly less than what could be extracted from the PS (Haarala et al. 1997) by treating the single-word definitions as synonyms, amounting to synonym candidates for 35067 out of the altogether 102740 entries, containing 506212 words, but compared to the PS the FINTHES database contains explicitly only synonyms, and its electronic version was/is considerably simpler to process as it has been supplemented with word-class data lacking from PS. Unfortunately, FINTHES has not been publicly documented. Another synonym dictionary of Finnish that must be mentioned is the one appended to NS, *Nyky-suomen sanakirjan synonyymisanakirja* (Jäppinen 1989), which contains some 18000 lexical entries. However, it appears that at least in the case of the studied THINK lexemes the synonym sets in this work correspond quite closely to the single-word definitions in PS, which is not that surprising as both works build upon NS. Therefore, I have relied on PS, even more so, as it was at my disposal in electronic form.

than exhaustive in nature.¹⁷ Nevertheless, even though the terminology and structure of Pajunen's general ontology and description has changed somewhat over time (cf. the tables/diagrams in Pajunen 1982: 336 in comparison to Pajunen 2001: 51-57, noted sketchily in Arppe 2006b), the conclusions as to the very close similarity of the argument structure of the selected THINK lexemes remain the same, though the picture has become more detailed.

In contrast, the second reference work (PS) has remained quite stable over the last decade, even more so as it is directed to a larger, non-professional audience. In its case, semantic similarity was assessed in terms of the extent to which the candidate synonyms shared the same words as definitions and the degree to which they could be judged substitutable with each other in the typically several usage examples given in the dictionary entries. In the end, this process had originally yielded several promising synonym groups, such as the THINK¹⁸ verbs *ajatella*, *mieltiä*, *pohtia*, *harkita*, and *tuumia/tuumata* 'think, ponder, consider, reflect' as well as the UNDERSTAND verbs *ymmärtää*, *käsittää*, *tajuta*, and *oivaltaa* 'understand, grasp, comprehend'.¹⁹

Out of these, the pair *mieltiä-pohtia* 'think, ponder' had been chosen (see Appendix A for an evaluation of the mutual interchangeability of these and the other THINK verbs in the example sentences given in the respective entries in PS), with their semantic similarity further validated by me through a manual assessment of the mutual interchangeability in each of the individual 855 sentences containing an instance of this verb pair in the originally used corpus (Keskisuomalainen 1994). The requirement satisfied by the *mieltiä-pohtia* pair was thus that of *strong entailment*, meaning that interchangeability applies for all (or practically all) cases. In this original selection process, it appeared to me that this strict criterion for the degree of substitutability will probably yield in larger numbers only *pairs* of synonyms, which are also common/frequent enough to be representative of general linguistic usage, and consequently the requirement would have to be relaxed somewhat in the case of synonym sets with more than two members. For instance, WordNet is based on a weaker notion of entailment, where interchangeability in at least some context(s) suffices for synonymy (Fellbaum 1998b: 77; Alonge et al. 1998: 21).

¹⁷ For instance, only *ajatella* and *mieltiä* are explicitly mentioned in Pajunen (1982), whereas *ajatella*, *tuumia*, *harkita*, *mieltiä*, and *järkeillä* are noted at various points in Pajunen (2001: 63, Table 8, 314, 317) (but still not *pohtia*).

¹⁸ Interestingly, THINK is one of the proposed *semantic primes* concerning mental predicates in *natural semantic metalanguage*, that is., *NSM*, (e.g., Goddard 2002, Table 1.2), a theory originally proposed by Wierzbicka (1996), although I was not aware of this at the time of the original selection process.

¹⁹ In the original frequency ranking using FINTHES and selecting only verb sets, the THINK verbs were to be found at rankings 51, 143, and 500 (with *mieltiä* and *pohtia* together at the last mentioned ranking); the UNDERSTAND verbs were to be found at ranking 217.

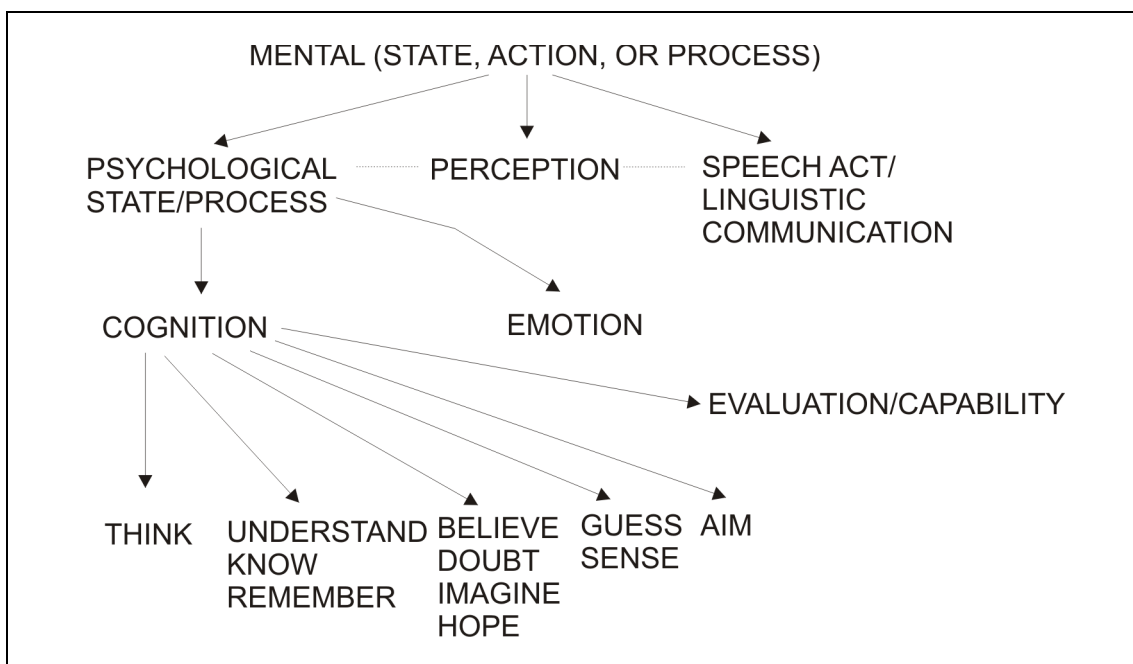


Figure 2.1. The semantic classification hierarchy of MENTAL verbs according to Pajunen (2001)

In order to establish and demarcate the extended group of synonymous THINK lexemes to be scrutinized in this dissertation, I will try to improve the above process by repeating it with a more comprehensive dictionary and a substantially larger corpus. I will first conduct a variant of the sub-test by using dictionary content from PS (Haarala et al. 1997), with the purpose of studying the overlap of word-definitions for lexemes belonging to the more general semantic grouping of COGNITION verbs, one step up from the THINK synonym group in Pajunen's (2001) hierarchy (see Figure 2.1). This will produce pair-by-pair ratings of similarity and dissimilarity for each lexeme against each other, for all those which are included in the analysis. Then, I will assess the resultant candidate synonym groupings with respect to the similarities/differences in the magnitudes of their relative frequencies, calculated on the basis of the largest corpus collection currently available for Finnish, namely, the Finnish Text Collection (FTC 2001).

2.1.2 Screening out the interrelationships of COGNITION lexemes by their dictionary definitions

I will begin with Pajunen's (2001: 313-314) treatment of COGNITION verbs, for which she distinguishes six different subclassifications. These consist of verbs of ANTICIPATION, e.g., *aanailta*, *ounastella* 'anticipate/foresee', *vaistota* 'sense', *uumoilla* 'guess [in a roundabout, unsure way]'; verbs of BELIEF and ASPIRATION, e.g., *epäillä* 'suspect/doubt', *haaveilla* 'dream', *kuvitella* 'imagine', *luulla* 'think/believe', *mieliä* 'aspire/desire to', *toivoa* 'wish/hope', and *uskoa* 'believe'; verbs of THOUGHT PROCESS and STATE OF KNOWLEDGE, e.g., *ajatella* 'think', *harkita* 'consider', *järkeillä* 'reason', *käsittää* 'grasp/understand', *mieltiä* 'think/ponder', *muistaa* 'remember', *tietää* 'know', *tuumia* 'think/reflect/muse', *ymmärtää* 'comprehend'; EVALUATIVE verbs, e.g., *arvioida* 'evaluate/judge', *huonoksua* 'consider bad', *väheksyä* 'belittle', *paheksua* 'consider improper', *paljoksua* 'consider as [too] much', *puntaroida*

‘weigh’, and *verrata* ‘compare’; verbs of INTENTION, e.g., *aikoa*, *meinata* ‘intend’, *suunnitella* ‘plan’, *tarkoittaa* ‘mean/intend for’; and verbs of ABILITY/CAPABILITY, e.g., *jaksaa* ‘have the strength to’, *kehdata* ‘dare/have the nerve to’, *kyetä* ‘can/have the capability to’, *onnistua* ‘succeed in’, *osata* ‘know how to’, and *pystyä* ‘can/be able to’. As the last subgroup is typically considered part of the Finnish modal verb system (see Kangasniemi 1992), I will exclude them from further scrutiny here.

Using these 31 (exemplified) COGNITION verbs explicitly mentioned by Pajunen (2001: 313-314) as a starting point, I first selected from PS all the dictionary entries in which any one of Pajunen’s examples was listed either as an entry lexeme or among the single-word definitions. This yielded 114 entries with 465 single-lexeme definitions, which consisted of 96 unique entries and altogether 168 unique lexemes, with which the same selection process was repeated once more, this time yielding 566 entries with 1498 single-lexeme definitions, representing 422 unique entry lexemes associated with 630 unique lexemes. If an entry lexeme had several explicitly indicated (i.e., numbered) distinct senses, that particular lexeme is listed repeatedly, each time together with those word-definitions that are associated with the sense in question. As a consequence, some word-definitions may also be counted in more than once for some particular entry lexemes, in such a case indicative of a shared range of senses with the recurrent word-definitions. An example of the thus extracted word definitions for our old friends *mieltiä* and *pohtia* is presented in Table 2.1. We can immediately see that *mieltiä* has slightly more individual definition words than *pohtia* (9 vs. 7); furthermore, as many as 6 are common for both, in addition to both named as a definition of the other. The full list of the selected COGNITION lexemes is given in Table B.1 in Appendix B, together with frequency information to be discussed below.

Table 2.1. Single-word definitions in PS for *mieltiä* and *pohtia*; common lexemes in boldface (no lexemes with repeated occurrences among the word-definitions).

Entry	Single-word definitions
mieltiä	punnita, harkita, ajatella, järkeillä, tuumia , mietiskellä, pohtia, suunnitella, aprikoida
pohtia	punnita, harkita, ajatella, järkeillä, tuumia , mieltiä, aprikoida

The purpose was to canvas in this manner any lexemes which in at least one of their senses could be used to denote a COGNITION state, process or activity. Therefore, no lexemes were excluded from the final set even though they obviously primarily denoted some other semantic field. Among these cases were, for instance, *nähdä* ‘see’ in the sense of UNDERSTAND, as if mentally “seeing”, *haistaa* ‘smell/sniff’, used figuratively as ‘get a whiff of something’, or *hautoa*, literally ‘incubate’ but also ‘hatch a plan (by oneself/in secret), foment, brood (long/alone)’. Furthermore, considering the entire COGNITION group instead of only the THOUGHT PROCESS subclass allowed also for assessing the degree of polysemy among the constituent lexemes, as some can clearly be considered to belong to more than one of the subclassifications, for instance *ajatella* as denoting both a THOUGHT PROCESS and INTENTION. My hypothesis was that quantitative analysis would link close in similarity those lexemes for which all the senses, or at least the primary ones, are associated primarily with COGNITION and any of its subclasses, while lexemes with multiple senses of which only one, possibly secondarily, concerns COGNITION, or which belong to more than one of its subclasses, would be relegated to the fringes.

With these word lists we can now quantify for each and every pairing of the selected COGNITION entry lexemes the extent of overlap among their definitions, which corresponds in principle to the sub-test (Church et al. 1994), *but* with single-word definitions used instead of significant collocates. The more word-definitions a lexeme has in common with another, the more intersubstitutable they can be considered, though this may be due to not only synonymy but also other types of lexical relationships between the two lexemes such as hyponymy or even antonymy, as Church et al. (1994) point out. However, using word-definitions instead of collocates, which Church et al. 1994 focused on, should specifically target synonymous lexemes. The resulting lists of lexemes similar in this respect with *mieltiä* and *pohtia* are presented in Table 2.2.

This time, we can see that *mieltiä* has common word-definitions with 25 of all the other COGNITION lexemes, whereas the corresponding figure for *pohtia* is slightly less at 23. Furthermore, in the case of both entries there are several other lexemes with which they share quite many word-definitions, indicating a closer relationship. Most notably, both share the most number of word-definitions with *ajatella*, *tuumia* and *aprikoida* in addition to each other. In fact, both *mieltiä* and *pohtia* have 19 lexemes (plus each other), with which they both share at least one word-definition, suggesting that the two lexemes would appear to be quite substitutable with one another. This overlap fits Pajunen's (62-63) assessment well that the classificatory structure of MENTAL verbs in general consists of lexical sets in which the members are in loose co-hyponymic relationships with each other. However, a substantial number of non-common lexical entries are also evident, which indicates that the lexemes are not exact synonyms in relation to each other.

Table 2.2. Overlap among the single-word definitions of *mieltiä* and *pohtia* with all the selected COGNITION lexemes; common lexemes in boldface.

Lexeme (number of lexemes with overlap)	Lexemes with overlap in definitions (number of overlapping items)
mieltiä (38)	ajatella (7), <i>pohtia</i> (6), tuumia (5), aprikoida (5), järkeillä (4), filosofoida (4), harkita (3), hautoa (3), funtsata (3), punnita (2), aikoa (2), tutkailla (2), tarkoittaa (2), tutkistella (2), spekuloida (2), meinata (2), meditoida (1), laatia (1), hankkia (1), tarkastella (1), ohjelmoida (1), katsoa (1), muistaa (1), pähkäillä (1), punoa (1), konstruoida (1), tuumailla (1), mitata (1), sommitella (1), arvella (1), mietiskellä (1), laskea (1), mitoittaa (1), tykätä (1), pohdiskella (1), keskustella (1), käsitellä (1), luonnostella (1)
pohtia (23)	ajatella (6), <i>mieltiä</i> (6), tuumia (4), aprikoida (4), funtsata (3), punnita (2), harkita (2), muistaa (2), järkeillä (2), filosofoida (2), hautoa (2), aikoa (1), katsoa (1), tuumailla (1), kelata (1), mitata (1), arvella (1), tarkoittaa (1), mietiskellä (1), laskea (1), spekuloida (1), tykätä (1), meinata (1)

In order to construct larger synonym sets, we could compare manually the overlap of the word-definitions for three, four or even more lexical entries. This is a feasible approach if we have a prior idea regarding which of the lexemes we want to consider (and thus also the size of the potential synonym set). In the case of the THINK lexemes, on the basis of my native speaker competence of Finnish and the two sets of overlapping lexemes presented in Table 2.2, I would be inclined to select as a

synonym set *ajatella*, *pohtia*, *tuumia*, *aprikoida*, *järkeillä*, *harkita*, *hautoa*, *punnita*, and *tuumailla*, with possibly also *filosofoida* and *funtsata*, the latter two lexemes being somewhat marked as sarcastic and slang terms, respectively. The word-definition overlaps for all of these entry lexemes, similar in form to those presented in Table 2.2, are presented in Table B.2 in Appendix B. In fact, this hypothesized synonym list overlaps with the synonym list anchored around *mieltä* (and shared by *ajatella*, *pohtia*, *harkita* and *tuumia*) in Jäppinen (1989), consisting, namely, of, ***ajatella***, ***mieltä***, ***mietiskellä***, ***pohtia***, ***pohdiskella***, ***harkita***, ***tuumia***, ***aprikoida***, ***järkeillä***, ***puntaroida***, ***punnita***, ***tuumata/tuumailla***, ***hautoa***, ***filosofoida***, ***meditoida***, ***spekuloida***, and ***funtsata/funtsia*** (where overlapping lexemes are in boldface). However, without such a hypothesis a blind exploratory comparison of all the possible permutations of trios and larger sets quickly becomes exceedingly large with even a relatively small number of lexical entries under overall consideration, with the number of permuted sets amounting to $n_{\text{permutations}} = n_{\text{entries}}! - (n_{\text{entries}} - n_{\text{set_size}})!$, and there would be no simple way to establish the proper size of synonym sets.

2.1.3 Clustering the COGNITION lexemes with statistical means

Under such circumstances, we may resort to a multivariate statistical method such as *hierarchical agglomerative cluster analysis (HAC)* (e.g., Kaufman and Rousseeuw 1990, see also Baayen 2007: 148-160), similar to what Divjak and Gries (2006) demonstrate, but by using either the single word definitions as such or the extent of their overlap as the classifying variables instead of contextual features derived from a corpus. A specific technique belonging to family of cluster analysis methods, HAC starts by considering all the items as singular clusters, which it then iteratively combines into larger clusters on the basis of maximizing intra-cluster similarity and minimizing inter-cluster similarity at each stage, ending up with a hierarchically nested tree structure. This data structure is typically represented as a so-called *dendrogram*, a sort of tree structure, which allows us to scrutinize visually the relationships of the individual items and then determine an appropriate set of clusters.

We can thus use this technique to cluster the entire set of selected COGNITION lexemes either according to 1) the single-word definitions as such and 2) the extent of overlap with respect to these single-word definitions, the complete results of which are presented in Figures B.1 and B.2 in Appendix B.²⁰ Interestingly, one can clearly discern in the overall dendrograms a distinct subcluster for THINK lexemes as well as another one for UNDERSTAND lexemes, with both of these sets being adjacent, and thus similar as groups, to each other. However, the overall hierarchy appears quite flat, and within the two subgroups “bushy”, which is in accordance with Pajunen (2001: 62-63, 313, see also Note 8 on page 434).

²⁰ As Divjak and Gries (2006: 37) note, there are several ways of calculating the similarity of the items and for determining the criteria for the amalgamation of the clusters, the selection of which significantly influences the resulting cluster structure. However, there are no deterministic, universally applicable rules for selecting these methods, which would guarantee an optimal solution. As Divjak and Gries have done, I have selected the methods that appear to produce the most useful results, these being the default *Euclidean* distance as a measure of (dis)similarity (in contrast to the *Canberra* method chosen by Divjak and Gries) and *Ward's rule* as the strategy for combining clusters (as did Divjak and Gries).

The two subclusters of THINK lexemes constructed with the two types of variables are presented in Figures 2.2 and 2.3, respectively. In general, we can see that the THINK subcluster based on the overlap corresponds exactly to the semantically hypothesized synonym set, whereas the subcluster based on the individual single-word definitions includes some additional lexemes. As these appear all in the overlap lists for *mieltiä* and *pohtia*, they can in some sense be used to denote the THINK concept. However, they are in my judgement either rarer and semantically quite specific lexemes, namely, *filosofoida* ‘philosophize, think philosophically/excessively theoretically (literally: to make philosophy out of something)’, *pohdiskella* ‘contemplate, ponder (aloud, now and then, not too seriously)’ or *pähkäillä* ‘think over (and over)’, or their primary sense is divergent from THINK “proper”, namely, *spekuloida* ‘speculate (out loud)’, *meditoida* ‘meditate’, *tutkailla* and *tutkiskella* ‘examine (study in one’s mind)’, *muistaa* ‘remember’, *meinata* and *tarkoittaa* ‘mean (intend to say)’. This judgement is also supported by their low degree of overlap. Interestingly, *ajatella* is clearly separated from all the rest in the overlap-based diagram (Figure 2.3), which could be explained by its role as the most frequent and prototypical of the group, as well as by its broad range of senses (see Section 2.3.2 below).

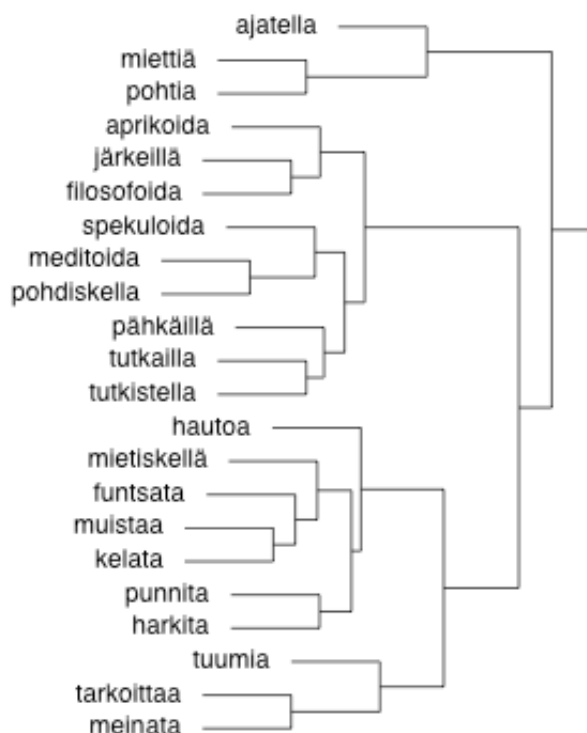


Figure 2.2. Subcluster of THINK lexemes on the basis of all the single-word definitions of the COGNITION lexemes.

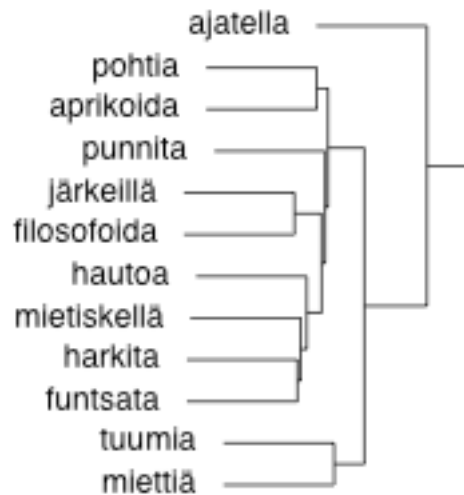


Figure 2.3. Subcluster of THINK lexemes on the basis of their overlap with respect to the single-word definitions.

2.1.4 Extracting frequencies for THINK lexemes

Next, I calculated the frequency rankings for all the verb lexemes in PS (Haarala et al. 1997) using the (base-form) lexeme frequency counts from the Finnish Text Collection (FTC 2001), the largest uniformly processed collection of Finnish to date. This corpus combines Finnish newspaper, magazine and literature texts from the 1990s, and amounts to some 180 million running text tokens. The corpus has been morpho-syntactically analyzed and disambiguated in its entirety using the *Textmorfo* parser (Jäppinen et al. 1983, Jäppinen and Ylilammi 1986, Valkonen et al. 1987) developed at Kielikone <URL: <http://www.kielikone.fi/>>; one should remember, however, that the results have not been manually verified and, what is more, no distinctions are obviously made between polysemous senses. In this corpus, a total of 25.4 million instances (roughly 14% of the running word tokens) were analyzed as verbs, representing 20930 distinct base-form lexemes. Of these, roughly over a half (12983) have at least two or more occurrences. This time, I used the natural logarithm of the relative frequency as an indicator for the magnitude of lexeme frequency (instead of the raw absolute or relative values as such), and the arithmetic average of these logarithm values as an indicator for the joint magnitude of the frequencies of a lexeme group as constituted by an entry and its single word-definitions (with only non-zero values included in the calculation). The frequencies and the individual and joint rankings of the 566 selected COGNITION entry lexemes and their associated single-word definitions are presented in full in Table B.1 in Appendix B.

We can now assess the relative frequencies of the subcluster of THINK lexemes identified above on the basis of the overlap in the single-word definitions (see Table 2.3). For the sake of comparison, Table 2.3 also contains the rankings from the Frequency Dictionary of Finnish (denoted FDF hereinafter) by Paunonen et al. (1979), which are also corpus-based figures.²¹ Interestingly, there is variation among

²¹ This Frequency Dictionary of Finnish is based on a corpus containing Finnish fictional texts, radio discussions, newspaper and magazine texts, and non-fiction reference works from the 1960s,

the individual rankings calculated here and those from the earlier source, though the orderings are overall quite similar: frequent lexemes in FTC are also frequent lexemes in the FDF, while infrequent lexemes here are again infrequent in FDF, if ranked at all.

As can be seen, the magnitudes as represented by the natural logarithms of the relative frequencies are very close for the three most frequent lexemes, namely, *pohtia*, *ajatella* and *miettiinä*, followed by *harkita* and *tuumia*, each alone on the next steps down on the magnitude ladder, before the rest of the more infrequent lexemes in the set. Indeed, visual inspection of Figure 2.4 also indicates that the observed frequencies of the scrutinized lexemes do not exactly conform to an ideal Zipfian distribution (see Appendix K), which is also supported by a goodness-of-fit test (with $P=0$).²² On the basis of these results, I decided to select for further study in this dissertation the four most frequent lexemes in the THINK group, namely, *ajatella*, *miettiinä*, *pohtia*, and *harkita*. In addition to clearly trailing *harkita*, the fifth-ranked *tuumia* has in comparison to the three most frequent lexemes only about one-tenth of occurrences. Furthermore, we will observe that in the final research corpus, described below in Section 2.4.2, *tuumia* has only 47 occurrences, which will be too low for the statistical analyses; the other THINK lexemes ranked as infrequent here have even fewer occurrences in the research corpus. Interestingly, it is exactly the selected set of four THINK lexemes which are given as the single-word definitions for the modern colloquial *funtsata/funtsia*. Moreover, we will later find out in Section 2.3.3 that three out of the four finally selected THINK lexemes, namely, *ajatella*, *pohtia* and *harkita*, have etymological origins in concrete activities of rural life particular to Finland, with the sole exception of *miettiinä* as a loan word. It will turn out that some vestiges of these concrete meanings can be interpreted to persist among the contextual preferences of the now abstract usages of these particular THINK lexemes.

amounting to slightly over 408 thousand words and representing 43670 base forms. Of these, 12663 (representing 90% of all the occurrences in the corpus) were selected for inclusion in the dictionary. Though this corpus is rather small by present standards, its selection of text types is quite representative, even more so as it contains a substantial amount of spoken language.

²² Furthermore, we may note that the ratio (0.371) of the most common (and assumedly semantically broadest) lexeme against all the rest is not exactly equal as Manin (submitted) has predicted for entire synonym groups, and in fact observed for the corresponding THINK lexemes in Russian, among others

Table 2.3. Absolute frequencies, the natural logarithms of the relative frequencies, and the corresponding ranking among verbs of the entire group of THINK lexemes identified on the basis of overlapping word-definitions in the PS, sorted according to descending frequency; ranks from FDF (Paunonen et al. 1979) include all word classes.

Lexeme	Absolute frequency	Natural logarithm of relative frequency	Ranking (among verbs)	Ranking in FDF (Paunonen et al. 1979)
pohtia	30572	-6.7	127	1792
ajatella	29877	-6.7	130	201
mieltä	27757	-6.8	141	1352
harkita	14704	-7.5	257	1063
tuumia	4157	-8.7	595	3740
punnita	2253	-9.3	828	3495
aprikoida	1293	-9.9	1153	11356
mietiskellä	995	-10.1	1345	9466
hautoa	536	-10.8	1939	4315
filosofoida	399	-11.1	2281	–
järkeillä	308	-11.3	2589	–
funtsata	29	-13.7	5996	–

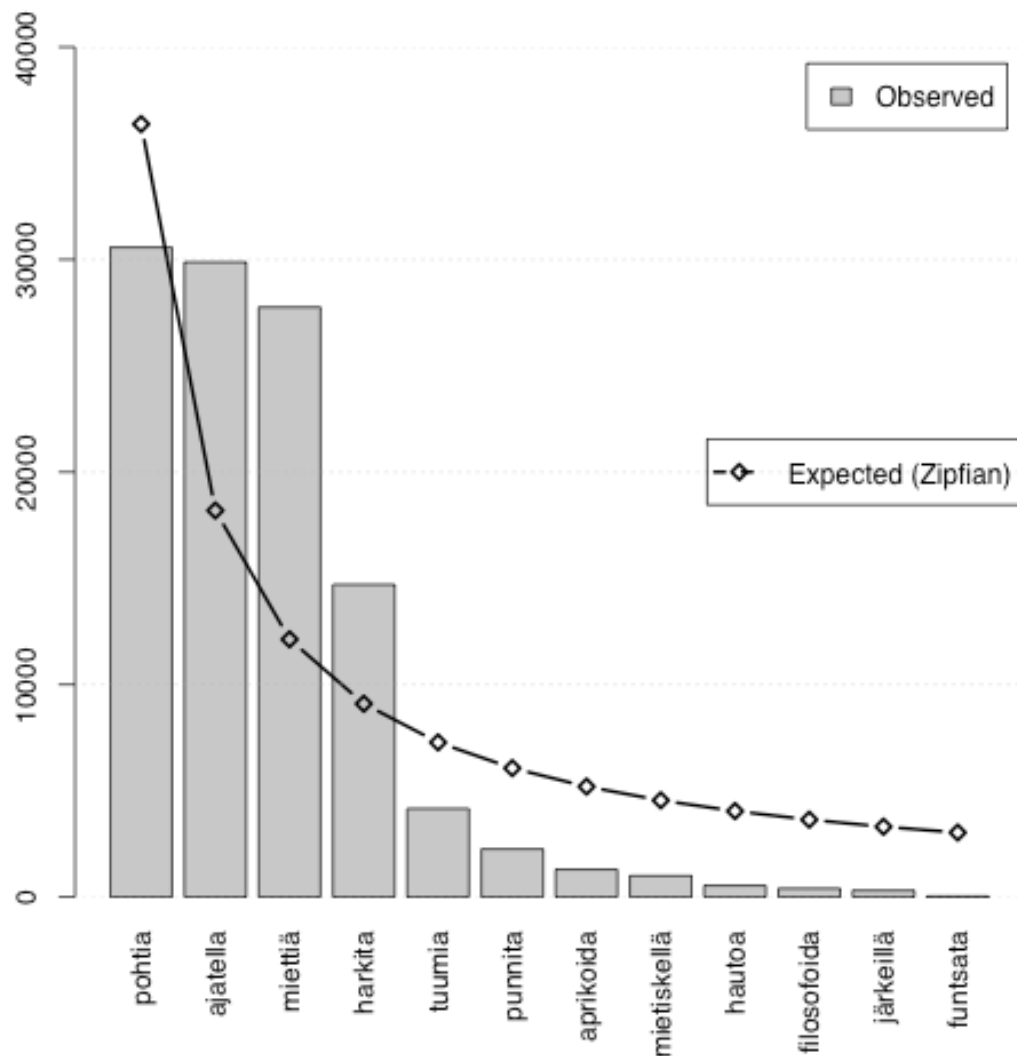


Figure 2.4. Frequencies of the entire group of THINK lexemes in FTC (2001), contrasted with an ideally Zipfian distribution of their joint frequencies.

2.2 Selection of the contextual features and their application in the analysis

The purpose in this study is to retrieve the entire contextual profile of the selected THINK lexemes, thus in principle following the Behavioral Profile approach advocated by Divjak and Gries (2006). However, I decided in practice to include all “verbal” uses, understood in the broad sense, for these lexemes, in contrast to Divjak and Gries (2006), who focus on one specific construction (FINITE form of the studied TRY verbs followed and modified by an INFINITIVE form of any verb). This covers, firstly, the FINITE (simplex) forms of the studied THINK lexemes, when they are used either alone or as FINITE or NON-FINITE auxiliaries in a verb chain. Secondly, this encompasses all the INFINITIVE and PARTICIPLE²³ forms of the studied THINK lexemes, including the

²³ Many structurally clearly participle forms in Finnish have usages when they can for all practical purposes be considered lexicalized adjectives, or to a lesser extent nouns, which becomes evident in translating them to, say, English. Often, they retain both a participial usage alongside the reading as an adjective or a noun, for instance, *kuollut* as both the adjective ‘dead’ and the participle on ‘die’ in

instances when these NON-FINITE verb forms are used as so-called *clause-equivalent constructions* (which in Finnish correspond to subordinate or relative clauses). However, nouns and adjectives derived from the studied verb lexemes, using morphemes traditionally considered as derivative in Finnish grammatical analysis, have been excluded from this study, though this is more of a formal delineation than a semantical one, for example, *ajatteleva* ‘thinking’ is included as a participle form while *ajattelematon* ‘unthinking/thoughtless’ is excluded as a derived adjective. Nevertheless, it would be perfectly possible to focus (later on) on some individual, specific constructional slot as Divjak and Gries (2006) have done.

The selection of contextual variables used in this study is rooted in traditional grammatical analysis, with the division into morphological, syntactic, semantic, pragmatic, discourse, and extra-linguistic levels, each with their own feature categories. Covering all of these feature category types, or at least as many as possible, within one and the same study is motivated by earlier research, which has firstly observed differences in the usage of synonyms within each category of features, and secondly interrelationships among feature categories at different levels, as was discussed earlier in Section 1.2. Out of the general list of possible analysis levels presented by Leech (1993, 2005), only the phonetic/prosodic and the pragmatic/discourse levels are clearly not included in this study. The omission of the first level follows quite naturally from the nature of the research corpus as consisting of only written text, albeit in many divergent modes; the lack of the latter level can be motivated as resulting from a focus on linguistic phenomena which are transparently apparent and observable in the linguistic structure of the text, without knowledge of the situational context (e.g., the attitudes or social relationships of the participants, which might be induced from a recording or personal participation). In spirit, this requirement of transparency and explicitness follows Divjak and Gries (2006: 35-36), although my set of features differs somewhat from the one they selected; this methodological affinity also holds for the consequent implicit restriction of the analyses to the immediate sentential context (Divjak and Gries 2006: 30). Moreover, such a broad selection of analysis levels and feature categories should conclusively address the critical formal analysis by Kenttä (2004) concerning results derived with the limited set of features in my earlier research concerning the THINK lexemes.

Furthermore, since my goal is not to demonstrate the suitability or superiority of one linguistic model or theory over another in describing the studied linguistic phenomenon, but rather to present a general methodology for combining a range of different feature categories from different levels of linguistic analysis in order to understand comprehensively the phenomenon in question, regardless of the underlying theory, I have generally opted for descriptive models which, on the one hand, have been recently applied to a range of languages including Finnish, and which have a computational implementation (if not yet for Finnish, then at least for some Standard Average European, i.e., ‘SAE’ language, as coined by Whorf 1956), on the other.

kuollut kieli ‘dead language’ vs. *hän on kuollut* ‘he has died/is dead’, or *tehtävä* as both the noun ‘task’ and a participle of the verb *tehdä* ‘do’ in *se on tehtävä* ‘it must be done’ vs. *onko mitään tehtävää?* ‘is there anything to be done/that should be done’ vs. *tehtävät asiat* ‘things to be done/that should be done’ vs. *vaikea tehtävä* ‘difficult task’ (vs. *se on vaikea tehtävä* ‘it is a difficult task/thing to do’ vs. *se on vaikea tehdä* ‘it is difficult to do’).

Indeed, such a stance of theory-neutrality (supported by, e.g., Leech 1993), or perhaps, a lack of passion for some particular theory, is facilitated by the fact that all of the recent major grammatical models, as reviewed in the overview by Hakulinen et al. (1994: 44-58) and thereafter, have been applied in one form or another for Finnish, albeit to differing degrees. With respect to a preference for computational resources, my underlying motivation has been to develop and test the methods presented in this dissertation which build upon tools, representations and resources that could later on be used to replicate similar analyses on a larger scale, once the resources in question attained sufficient quality (being parsing in the case of Finnish) or become localized also for Finnish (being semantic ontologies of the WordNet type). This general idea of developing linguistic theory in close interaction with its computational implementation and empirical performance can be considered a defining characteristic of the “Helsinki school” of linguistic analysis in the 1980-90s, exhibited in morphology (and also phonology) by the *Two-Level model* (TWOL) by Koskeniemi (1983), and in syntax first by the *Constraint Grammar* (CG) formalism by Karlsson et al. (1990, 1995), and later the related *Functional Dependence Grammar* (FDG) formalism by Tapanainen and Järvinen (1997, 1998); and it is an attitude to which I, too, subscribe and attempt to follow in this study.²⁴

The selection of a linguistic model and the categorizations and features it incorporates, and their subsequent application in the linguistic analysis, that is, *annotation* of the data, are closely intertwined. As Leech (1993) notes, annotation is not an absolute but an *interpretative* activity. Only when one applies the originally selected model and features to the real data at hand does one learn how well they suit the studied phenomenon and are able to describe it, and what unanticipated new aspects evident in the data one would also like to or need to cover. Linguistic categories and features are typically defined by prototypical usage cases; if one wants to keep to the original feature set, and the generality it represents, rather than create novel features to cater to new circumstances and special cases, one would have to creatively modify, bend, and extend the original definitions, while at the same time maintaining the validity of the prior classifications. Sampson (1995: 14-16) has aptly described this process in my view, using an analogy with the legal system of *common law*, which can be characterized by the constant consideration of the applicability, or lack thereof, of precedents and of the possible need to set new ones. However, one will soon notice the emergence of regularities and a sort of convergence in that only a subset of all possible features and their combinations account for most of the context for any set of semantically similar words scrutinized at a time, which also Divjak and Gries (2006) note. Though all theoretically possible combinations of features representing possible contexts might at first glance appear close to infinite, in practice the number of actually observable context types is quite finite and manageable. This was the case with the linguistic analysis of the context of the selected THINK lexemes in this study.

The details of the various stages and levels of linguistic analysis applied to the research corpus are covered at length in Appendix C, but I will briefly cover the main points also here. The research corpus was first automatically morphologically and

²⁴ This approach is in fact quite similar to the one adopted by Divjak and Gries (2006), who, in their quantitative study of Russian synonymy, implement (a substantial part of) the comprehensive descriptive methodology developed and applied for Russian lexicography by the Moscow School of Semantics (Apresjan and others).

syntactically analyzed using a computational implementation of Functional Dependency Grammar (Tapanainen and Järvinen, 1997, Järvinen and Tapanainen 1998) for Finnish, namely, the FI-FDG parser (Connexor 2007). Thus, the morphological analysis employed in this study can be characterized as compositional, based on traditionally-defined atomistic morphological features, and the syntactic analysis as monostratal, based on directly observable surface structure and consisting of dependency relationships between elements representing various functional roles. Moreover, such syntactic elements may consist of multiple words and can be discontinuous.

After the automatic analysis, all the instances of the studied THINK lexemes together with their syntactic arguments were manually validated and corrected, if necessary, and subsequently supplemented with semantic classifications by hand. Each nominal argument (in practice nouns or pronouns) was semantically classified into one of the 25 top-level *unique beginners* for (originally English) nouns in WordNet (Miller 1990). Furthermore, subordinate clauses or other phrasal structures assigned to the PATIENT argument slot were classified following Pajunen (2001) into the traditional types of participles, infinitives, indirect questions, clause propositions indicated with the subordinate conjunction *että* ‘that’ and direct quotes with attributions of the speaker using one of the studied THINK lexemes (e.g., “...” *mieltii/pohtii joku* “...” thinks/ponders somebody’). This covered satisfactorily AGENTS, PATIENTS, SOURCES, GOALS, and LOCATIONS among the frequent argument types as well as INSTRUMENTS and VOCATIVES among the less frequent ones.

However, other argument types, which were also frequent in the context of the studied THINK lexemes indicating MANNER, TIME (as a moment or period), DURATION, FREQUENCY, and QUANTITY, consisted of a high proportion of adverbs, prepositional/postpositional phrases, and subordinate clauses (or their CLAUSE-EQUIVALENTS based on NON-FINITE verb forms). These argument types were semantically classified following the *ad hoc* evidence-driven procedure proposed by Hanks (1996), in which one scrutinizes and groups the individual observed argument lexemes or phrases in a piece-meal fashion, as the contextual examples accumulate, and thus generalizes semantic classes out of them, without attempting to apply some prior theoretical model. Only in the case of MANNER arguments did there emerge several levels of granularity at this stage in the semantic analysis. Moreover, even though clause-adverbials (i.e., META-comments such as *myös* ‘also’, *kuitenkin* ‘nevertheless/however’ and *ehkä* ‘maybe’, as well as subordinate clauses with *mutta* ‘but’ and *vaikka* ‘although’) were also relatively quite frequent as an argument type, they were excluded at this stage due to their generally parenthetical nature.

Furthermore, as an extension to Arppe (2006b), the verb chains, of which the studied THINK lexemes form part, were semantically classified with respect to their modality and other related characteristics, following Kangasniemi (1992) and Flint (1980). Likewise, those other verbs which are syntactically in a co-ordinated (and similar) position in relation to the studied THINK lexemes were also semantically classified, following Pajunen (2001). Moreover, with respect to morphological variables, I chose to supplement analytic features characterizing the entire verb chain of which the studied THINK lexemes were components, concerning polarity (i.e., NEGATION vs. AFFIRMATION), voice, mood, tense and person/number. In addition, in a further abstraction in comparison to Arppe (2006b), the six distinct person/number features

(e.g., FIRST PERSON SINGULAR, FIRST PERSON PLURAL, SECOND PERSON SINGULAR, and so on) were decomposed as a matrix of three person features (FIRST vs. SECOND vs. THIRD) and two number features (SINGULAR vs. PLURAL). Finally, with respect to the extra-linguistic features, these concerned the two sources, representing to distinct media, which the research corpus consisted of, their constituent subdivisions, author designations, and various aspects of repetition of the selected THINK lexemes within individual texts.

2.3 Present descriptions of the studied THINK synonyms

We may now turn to the specific descriptions of the studied THINK lexemes in the external reference sources presented in Section 2.1.1, namely, the general description of the Finnish verb lexicon by Pajunen (2001), *Argumenttirakenne* ‘Argument structure’, as well as their lexical entries in two current dictionaries, *Perussanakirja* (Haarala et al. 1994-1997, Haarala et al. 1997) ‘Standard Dictionary of Finnish’ and *Nykysuomen sanakirja* (Sadaniemi et al. [1951-1961] 1976) ‘Dictionary of Modern Finnish’. At times, I will refer to these in the following discussion by their acronyms *AR*, *PS*, and *NS*, respectively.

2.3.1 THINK lexemes in Pajunen’s *Argumenttirakenne*

According to Pajunen (2001: 313-319), COGNITION verbs, under which the THINK lexemes belong in her classificational hierarchy, have typically two arguments, but in conjunction with comparative or evaluative readings they may take a third, additional argument. Table 2.4 below presents the so-called *lexicon forms* that in *AR* formally represent the argument structure and that are particular to the studied THINK lexemes. In the first place, we can note that Pajunen distinguishes *harkita* ‘consider’ from the rest in terms of its argument context, while she considers the most prototypical of the group, *ajatella* ‘think’, as similar to *käsittää* ‘understand’ in these respects. However, she does not explicitly state under which of the presented lexicon forms the other THINK lexemes considered in this study, that is, *mieltiä* and *pohtia*, should fall.²⁵ Nevertheless, the actual differences between these presented two lexicon forms are not that significant: for both, the first argument (X-ARG) corresponds to the syntactic subject and the second argument (Y-ARG) to either the syntactic object or a clause argument (*lausemäärite*, abbreviated in Pajunen’s notation as *LM*). However, the agentivity typical to the first argument is slightly weaker for *harkita* than for *ajatella* (and *käsittää*, for that matter). Thus, volitional participation in a (mental) state or event is stronger for *ajatella* than *harkita*, while sensing and/or perceiving is equally characteristic of both. Nevertheless, in the overall perspective the agentivity of the COGNITION verbs, and consequently also of the THINK lexemes as its subgroup, is quite weak (Pajunen 2001: 300).

Furthermore, the range of clause arguments as the second argument (Y-ARG) which is possible for *ajatella* and its kind is somewhat broader than that for *harkita*. While both can in this position instead of a syntactic object also take a subordinate clause and a PARTICIPIAL construction, which may have either a common or a disjoint subject with the first argument (X-ARG) (though the corpus-based observations of *harkita* exhibit a categorical preference for joint subjects, Pajunen 2001: 405-406, Table 41), only *ajatella* can have an INFINITIVE in this argument slot. This characteristic associates *ajatella* with the third lexicon form shown in Table 2.4, namely, pertaining to *aikoa* ‘intend’, and this sense is also apparent among the dictionary entries presented below for *ajatella* as well as *mieltiä*. Moreover, while the events or states denoted by the PARTICIPIAL constructions in the second argument position for *harkita*

²⁵ My linguistic intuition as a native speaker of Finnish would place *mieltiä* together in the same lexicon form as *ajatella*, and *pohtia* instead in the same lexicon form with *harkita*, in that a (first) infinitive as the second argument (Y-ARG) would in my judgment seem rather odd with *pohtia*, but to at some extent conceivable with *mieltiä*.

must be asynchronous (in this case temporally posterior) with the main verb, that is, *harkitsin lähteväni* ‘I considered leaving [some time following the consideration]’ vs. *harkitsin *tulleeeni* ‘I considered to have come’, this is not obligatory for *ajatella*, that is, *ajattelin lähteväni* ‘I thought of leaving [intended to leave]’ vs. *ajattelin tulleeeni* ‘I thought of having come’ (Pajunen 2001: 405, though her later corpus-based observations for *ajatella* exhibit >90% preference for synchronous PARTICIPIAL constructions, Pajunen, 2001: 407, Table 42).

Table 2.4. Adaptation of the *lexicon forms* for the studied THINK lexemes in Pajunen (2001: 316-318, specifically Table 48, page 317).

käsittää, ajatella ‘understand, think’:

X-ARG: Subject, Agentivity: volitional participation in state or event; sensing and/or perceiving

Y-ARG: Object, Clause Argument=subordinate clause, participial construction (common or disjoint subject with X-ARG²⁶), infinitive

harkita ‘consider’:

X-ARG: Subject, Agentivity: (volitional participation in state or event); sensing and/or perceiving

Y-ARG: Object, Clause Argument, participial construction (disjoint or common subject²⁷ with X-ARG, possibly asynchronous)

voida, aikoa ‘can/may, intend’:

X-ARG: Subject

Y-ARG: Clause Argument: Infinitive

In general, among Finnish verbs, COGNITION lexemes appear to have the broadest range in the types of clause arguments that they can take as their second argument (Y-ARG), including in practice all the available types of syntactic structures, namely, INFINITIVES, PARTICIPIAL constructions, INDIRECT QUESTIONS, and *että* ‘that’ clauses, the latter two being both subordinate clauses (Pajunen 2001: 358-360). As I have noted in Appendix C concerning the use of THINK verbs as attributive phrases in conjunction with citations, similar to SPEECH ACT verbs, we can in my view also include independent clauses among their acceptable clause arguments, at least structurally speaking. This view is in accordance with cross-linguistically derived syntactic frames available for THINK, when considered as a semantic prime following the Natural Semantic Metalanguage (NSM) approach (Goddard 2003: 112, Example 1c) to be discussed further below, as well as with Fortescue’s (2001: 28-30) observations that THINK lexemes in many languages have also developed a polysemy meaning ‘say/pronounce’, or that they may have originated from such words by metaphorical abstraction. However, this is in some contrast to Pajunen (2001: 363-366, 428-430) who rules such usage out on the grounds that it renders COGNITION verbs semantically as parenthetical expressions (cf. also Atkins and Levin 1995: 106-107 concerning an even broader class of lexemes, for example, *sniff*, *snort*, *bellow*,

²⁶ Pajunen’s (2001: 407, Table 42) own corpus data would suggest that participial constructions as the second argument (Y-ARG) for *ajatella* exhibit a strong preference (>90%) for synchronicity with the node verb, which she has nevertheless omitted from this lexicon form.

²⁷ Pajunen’s (2001: 317, Lexical form 48) judgment of preference order appears to be in contradiction with the actual preferences observed in her corpus (Pajunen 2001: 406, Table 41).

murmur as well as *shudder* and *quiver*, all concomitant with [some manner of] speech, which may be used similarly in an attributive way). A key aspect in the different types of clause arguments is that they vary in the extent to which they can indicate tense and mood in general and in relation to the verb of which they are an argument, so that the independent clauses as well as the various subordinate clauses can mark both tense and mood, PARTICIPIAL constructions only tense (with limitations), while INFINITIVES (as well as deverbal nominal derived forms) are entirely bare in this respect.

With regard to the semantic types of the arguments of the COGNITION lexemes, overviewed in Table 2.5, the first arguments (i.e., syntactic subjects) are, according to Pajunen (2001: 316-318), without exception HUMAN referents. In turn, the second argument has greater variety, denoting typically a CONCRETE OBJECT or an ABSTRACT NOTION or STATE-OF-AFFAIRS, and sometimes also ANIMATE ENTITIES. As a functional PATIENT in general, the second argument may alternatively refer to thought(s) stimulated by the external reality, having thus directionality from the world into the mind, or the result of cognitive activity, where the directionality is reversed to flowing from the mind to the world. Specifically with respect to syntactic objects as the second argument, these refer, according to Pajunen (2001: 316-317), mostly to ABSTRACT NOTIONS, while CONCRETE or ANIMATE referents are not fully applicable here with all COGNITION lexemes, and the use of HUMAN referents are natural only with a subset of the lexemes, and furthermore often in only restricted contexts. Concerning the third possible argument (Z-ARG), topic or (discourse) subject referents are mentioned as possible by Pajunen (2001: 318), that is, *ajatella jotakin jostakin asiasta* ‘think something about/concerning some matter’, as well as evaluative or comparative statements of various kinds, that is, *ajattelin hänen ymmärtävän asian* ‘I thought him to understand the matter’, or *ajattelin häntä viisaammaksi [kuin hän oli]* ‘I thought him wiser [than he was]’.

Table 2.5. Semantic classifications associated with the syntactic arguments in the *lexicon forms* for THINK lexemes presented above in Table 2.4, following Pajunen (2001: 316-318)

X-ARG: \forall human referent
Y-ARG: concrete entity, abstract notion, state-of-affairs > animate entity as referent; stimulus (‘world-to-mind’); result (‘mind-to-world’); Object: abstract notion > concrete object, state-of-affairs, animate/human referent
Z-ARG: subject/topic also in conjunction with comparative or evaluative usage (translative construction)

We may in conclusion note that Pajunen (2001) does not address the potential morphological preferences of the verbs at all (of the type observed by Arppe [2002] regarding the studied THINK lexemes). What is more, she suggests no differentiation among the various types of human referents, such as have been observed by Arppe and Järvikivi (2002, 2007b). Finally, characteristic associations with possible syntactic arguments other than the three basic types covered above (i.e., the obligatory X-ARG and Y-ARG, as well as the optional Z-ARG) are not asserted.

At this point, we can also compare Pajunen’s description, a study particular to Finnish however language-typologically oriented it strives to be, against the cross-linguistic conclusions derived within the natural semantic metalanguage (NSM) approach (e.g.,

Goddard 2002; Wierzbicka 1996). In this framework, four syntactic frames are considered to be universally available to the semantic prime THINK (Goddard 2003: 112), presented in 2.1 below with Pajunen’s syntactic argument types supplemented, when possible. Of these frames, (2.1a) and (2.1d) can be seen to correspond with Pajunen’s two-slot lexicon forms (consisting of X-ARG and Y-ARG), with either a (nominal) object or an *että* ‘that’ subordinate clause as the second argument Y-ARG, while (2.1c) extends this set of clause argument types to include entire clauses, as was discussed above. In view of the syntactic structures not expressly dealt with by Goddard (2003), if one considers PARTICIPIAL constructions used as CLAUSE-EQUIVALENTS to be equivalent to *että*-clauses, they could be placed under frame (2.1d); however, INDIRECT QUESTIONS do not appear to have an obviously natural home among these frames. Furthermore, frame (2.1b) would conform to Pajunen’s three-slot lexicon form (consisting of X-ARG, Y-ARG and Z-ARG) associated with evaluative and comparative statements.

- (2.1) a. X_{X-ARG} thinks about Y_{Y-ARG} [topic of thought]
 b. X_{X-ARG} thinks something_{Y-ARG} (good/bad) about Y_{Z-ARG} [complement]
 c. X_{X-ARG} thinks like this: “_Y-ARG” [quasi-quotational complement]
 d. X_{X-ARG} thinks {that []_S}_{Y-ARG} [propositional complement]

2.3.2 THINK lexemes in *Suomen kielen perussanakirja* and *Nykysuomen sanakirja*

Next, we may move on to see what inter-lexical semantic and contextual syntactic information the two current Finnish dictionaries contain with respect to the studied THINK lexemes. Both dictionaries contain four types of information for each lexical entry, exemplified in Table 2.6 for *pohtia* as defined and described in *Perussanakirja* (Haarala et al. 1997). In conjunction with the head word (field 1), which in the case of verbs is traditionally presented in the FIRST INFINITIVE form in Finnish dictionaries, we can find a code (field 2) indicating the inflectional (verbal) paradigm to which *pohtia* belongs, thus being similar to *lähteä* ‘leave’ and having the consonant gradation alternation *F: t ~ d*, for example, *pohtia* ‘[to] think’ vs. *pohdin* ‘I think’. This is followed by the definition proper (field 3), consisting to a large part of singular words which have at least one sense in common with the lexical entry²⁸, but often also initiated by a multiword qualification constructed around some more general, prototypical word representing the same semantic field, being in this case *ajatella*. This would fit perfectly within the NSM framework in which THINK is one of the (universal) semantic primes around which other words are defined (Goddard 2003). In these terms, *pohtia* is principally defined quite elaborately as *ajatella jotakin perusteellisesti, eri mahdollisuuksia arvioiden* ‘think about something thoroughly, evaluating different possibilities’. The fourth and last field in the lexical entry provides one or more example sentences or fragments, which in the case of PS are currently corpus-based. In contrast, the examples in NS, dating from the middle of the 20th century, have been selected from a vast collection of citation cards, often representing idiomatic usage by established Finnish authors such as Aleksis Kivi, F. E. Sillanpää or Volter Kilpi, or otherwise commonly known works such as the *Kalevala* or the *Bible* (Sadaniemi [1965] 1976: vi). In both dictionaries, the examples are quite often constructed around the canonical FIRST INFINITIVE, being thus

²⁸ These are possibly but not necessarily synonyms of the head word (Haarala et al. 2000: xxi).

practically AGENT-less, which is the case with the three example fragments in Table 2.6; likewise, the subjects in fragments with FINITE verb forms are almost always omitted, as the AGENT is manifested in the inflectional form, although this is in normal language usage in principle correct only in the case FIRST and SECOND PERSON and impersonal THIRD PERSON SINGULAR forms. Finally, if the head word of a lexical entry is associated with more than one distinct sense, each of these has its own definition(s) and example(s), but that is not the case here with *pohtia*.

However, no explicit syntactic information about the possible argument contexts is provided in PS nor in NS, in contrast to, for example, the *Collins COBUILD English Dictionary* for English (Sinclair et al. 2001), nor are possible morphological preferences discussed. As was already noted above in Section 2.1.1, PS is in many respects a revised and updated version of NS, and thus both dictionaries can clearly be observed to share a great deal in terms of their definitions and usage examples for lexical entries, while the differences between the two dictionaries are mostly due to changes in the Finnish language, culture, and society in the last 50 years (characterized by the transition from a predominantly agrarian and rural nation to an urban industrialized one), as well as to the more concise selection of lexical entries in *Perussanakirja* (being roughly half that of NS, see Haarala et al. 1990: v-vi).

Table 2.6. Original Finnish lexical entry for *pohtia* in *Perussanakirja* (PS), with the four component fields marked out, and the multiword definition underlined.

<p>[1/LEXICAL ENTRY: pohtia] [2/INFLECTIONAL PARADIGM CODE: 61*F] [3/DEFINITION: <u>ajatella jotakin perusteellisesti, eri mahdollisuuksia arvioiden, harkita, miettiä, tuumia, ajatella, järkeillä, punnita, aprikoida.</u>] [4/USAGE EXAMPLES: <i>Pohtia arvoitusta ongelmaa. Pohtia kysymystä joka puolelta. Pohtia keinoja asian auttamiseksi.</i>]</p>

We can primarily use the definitions provided in the lexical entries to sketch out the meaning potentials and similarity associations of the studied THINK lexemes, a concise outline of which is presented in Table 2.7 (with approximate English translations immediately below in Table 2.8), based on the more current of the two sources, namely, PS. As can be seen, *ajatella* has by far the largest number of senses (5 or 6, depending on whether one includes the specialized exclamative usage or not), while *harkita* and *miettiä* have two each and *pohtia* only one; nevertheless, the primary meaning for each of the selected four THINK lexemes is defined using all the other three, among others. Furthermore, there are several less frequent THINK lexemes which are shared as definitions among the selected four. Among these, *punnita* ‘weigh’ is common for all, while *tuumia*, *aprikoida* ‘think’, and *järkeillä* ‘reason’ are shared by all but one in various constellations. In addition, *ajatella*, *miettiä*, and *harkita* have in common as a secondary sense *suunnitella* ‘plan/intend’, which has been observed as a common metonymic extension in a range of languages for THINK lexemes (Fortescue 2001: 26-26, 38, Goddard 2003: 116); however, *pohtia* lacks this future-oriented characteristic.

In addition to these shared senses, *ajatella* can also be used to denote having or constructing an opinion or attitude concerning something, thus corresponding meaning-wise to *asennoitua*, *suhtautua* and *arvella*, imagining, assuming/presuming or presupposing something, associated then with *kuvitella*, *olettaa*, and *otaksua*, which generally speaking all fit Fortescue’s (2001: 28) cross-linguistic observations

of ‘believe’ as a common polysemous extension of the THINK lexemes, or a more focused or long-term direction of cognitive attention towards some concrete or abstract entity. Although the other three THINK lexemes do have some of these aforementioned qualities, they could not replace *ajatella* in the associated usage examples in my judgement as a native speaker of Finnish. Moreover, *harkita* has a further secondary sense denoting reaching or ending up with a (mental) conclusion through thorough consideration, which may at first glance seem distinct from the rest; however, this can be understood simply as a more conscious, objective, and drawn-out form of the ‘opine’ sense of *ajatella*.

Table 2.7. Original Finnish definitions of the studied THINK lexemes in *Perussanakirja* (PS); single-word definitions common to at least three underlined, those common to all four marked in addition in **boldface**.

<p>ajatella (5-6)</p> <p>1. yhdistää käsitteitä ja mielteitä tietoisesti toisiinsa (usein jonkin ongelman ratkaisemiseksi), <u>mieltiä</u>, <u>harkita</u>, <u>pohtia</u>, <u>tuumia</u>, <u>järkeillä</u>, päätellä, <u>aprikoida</u>, punnita.</p> <p>2. asennoitua, suhtautua, olla jotakin mieltä jostakin, arvella.</p> <p>3. kuvitella, olettaa, pitää mahdollisena, otaksua.</p> <p>4. kiinnittää huomiota johonkin, ottaa jotakin huomioon, pitää jotakin silmällä, mielessä.</p> <p>5. <u>harkita</u>, aikoa, <u>suunnitella</u>, <u>tuumia</u>.</p> <p>6. vars. ark. huudahduksissa huomiota kiinnittämässä tai sanontaa tehostamassa.</p>	<p>mieltiä (2)</p> <p>1. <u>ajatella</u>, <u>harkita</u>, <u>pohtia</u>, punnita, <u>tuumia</u>, <u>aprikoida</u>, <u>järkeillä</u>, mietiskellä.</p> <p>2. <u>suunnitella</u>; keksiä (miettimällä).</p>
<p>harkita (2)</p> <p>1. <u>ajatella</u> perusteellisesti, eri mahdollisuuksia arvioiden, <u>pohtia</u>, punnita, <u>puntaroida</u>, <u>mieltiä</u>; <u>suunnitella</u>.</p> <p>2. päätyä johonkin perusteellisen ajattelun nojalla, tulla johonkin päätelmään, katsoa joksikin.</p>	<p>pohtia (1)</p> <p><u>ajatella</u> jotakin perusteellisesti, eri mahdollisuuksia arvioiden, <u>harkita</u>, <u>mieltiä</u>, <u>tuumia</u>, <u>ajatella</u>, <u>järkeillä</u>, punnita, <u>aprikoida</u>.</p>

Table 2.8. Approximate English translations for the definitions of the studied THINK lexemes in *Perussanakirja* (PS); single-word definitions common to at least three underlined, those common to all four marked in addition in **boldface**.

<p>ajatella (5-6)</p> <p>1. think/contemplate/reflect, consider/deliberate, <u>ponder</u>, <u>deem</u>, <u>reason</u>, deduce, <u>riddle</u>, weigh</p> <p>2. regard, relate to, have some opinion concerning something, suppose/believe/guess</p> <p>3. Imagine, assume/presume, consider possible, presuppose.</p> <p>4. Focus attention on something, take something into consideration, keep an eye on something, keep something in mind.</p> <p>5. <u>Consider</u>, intend, <u>plan</u>, deem.</p> <p>6. [Colloquial: In exclamations to attract attention or intensify the expression].</p>	<p>mieltiä (2)</p> <p>1. <u>think</u>, <u>consider</u>, <u>ponder</u>, weigh, <u>deem</u>, <u>riddle</u>, <u>reason</u>, meditate.</p> <p>2. <u>plan</u>; conceive of (by thinking).</p>
<p>harkita (2)</p> <p>1. <u>think</u> thoroughly, evaluating different alternatives/possibilities, <u>ponder</u>, weigh, [weigh], [think]; <u>plan</u>.</p> <p>2. conclude something on the basis of thorough thinking, end up with some conclusion, consider as something.</p>	<p>pohtia (1)</p> <p>1. <u>think</u> about something thoroughly, evaluating different possibilities, <u>consider</u>, [think], <u>deem</u>, think, <u>reason</u>, weigh, <u>riddle</u>.</p>

Secondly, we can scrutinize the usage examples in order to see what information they implicitly encode with respect to syntactic and semantic contextual preferences of each of the studied THINK lexemes. In practice, this amounts to treating the usage examples as if they constituted a very representative sample, that is, a concise corpus, concerning the studied THINK lexemes, which is certainly what one could expect of a (corpus-based) dictionary. Tables 2.9 and 2.10 represent the linguistic analyses of the original Finnish lexical entry for *pohtia* in both PS and NS, using the array of contextual feature variables presented in depth in Appendix C; approximate English translations of these analyses are provided in Tables 2.11 and 2.12. Corresponding treatments for the other three THINK lexemes, namely, *ajatella*, *mieltiä* and *harkita*, are presented in Appendix F. In addition to the actual example sentences and fragments, arguments in multiword definitions have also been analyzed.

Firstly, we can see that the lexical entry in NS contains as the first sense for *pohtia* the original agrarian meaning ‘winnow’, which is no longer present in the more modern PS. Nevertheless, we should make a mental note at this point of the PASSIVE voice exhibited in the singular usage example for this older sense, that is, *Vilja pohdittiin_{PASSIVE} pohtimella* ‘The grain was winnowed_{PASSIVE} with a winnower’. Secondly, we can see among the examples for the more abstract (and currently more common) sense of *pohtia* represented in both dictionaries one shared example fragment (underlined), demonstrating the continuity between the two dictionaries. We may also note that for this second sense NS has more and longer examples than PS (5 vs. 3); furthermore, three of these in NS are complete sentences, in comparison to PS where all examples are clausal fragments constructed around the canonical FIRST INFINITIVE form. In all, we can observe quite an amount of contextual information, with 3 occurrences of 1 unique morphological feature and 9 occurrences of 5 distinct couplings of syntactic arguments and their semantic classifications among PS’s

examples and definitions. In NS, the respective figures (excluding the older agrarian sense) are somewhat higher, with 13 occurrences of 9 unique morphological features and 11 occurrences of 8 distinct argument-classification couplings, and 3 occurrences of semantically un-classified arguments.

Table 2.9. Original lexical entry in Finnish for *pohtia* in *Perussanakirja* (PS) and its linguistic analysis; default lexical entry forms (i.e., sentence-initial FIRST INFINITIVES) in parentheses; examples common with NS underlined.

pohtia^{61*F}

ajatella jotakin_{PATIENT+NOTION?} perusteellisesti_{MANNER+THOROUGH}, [eri mahdollisuuksia arvioiden]_{MANNER+THOROUGH}, harkita, miettiä, tuumia, ajatella, järkeillä, punnita, aprikoida.

*Pohtia*_(INFINITIVE1) arvoitusta_{PATIENT+NOTION/COMMUNICATION} ongelmaa_{PATIENT+NOTION}.

*Pohtia*_(INFINITIVE1) kysymystä_{PATIENT+COMMUNICATION} joka puolelta_{MANNER+THOROUGH}.

Pohtia_(INFINITIVE1) keinoja_{PATIENT+ACTIVITY/(NOTION)} asian auttamiseksi_{PURPOSE/REASON+ACTIVITY}.

Table 2.10. Original lexical entry in Finnish for *pohtia* in *Nykysuomen sanakirja* (NS) and its linguistic analysis; default lexical entry forms (i.e., sentence-initial FIRST INFINITIVES) as well as default features (i.e., ACTIVE voice and SINGULAR number) in parentheses; examples common with PS underlined.

pohtia^{17*} (verbi)

1. (=pohtaa) | *Vilja*_{PATIENT+SUBSTANCE} pohdittiin_{ANL_PASSIVE+ANL_INDICATIVE+ANL_PAST} pohdimella_{INSTRUMENT+ARTIFACT}. -- (tavallisesti) 2. harkita, miettiä, tuumia, ajatella, järkeillä, punnita, aprikoida | *Pohtia*_(INFINITIVE1) jotakin seikkaa_{PATIENT+NOTION}, tilannetta_{PATIENT+STATE}.
*Pohtia*_(INFINITIVE1) keinoja_{PATIENT+ACTIVITY/(NOTION)} jonkin asian auttamiseksi_{PURPOSE/REASON+ACTIVITY}.
*Kysymystä*_{PATIENT+NOTION/COMMUNICATION} pohdittiin_{ANL_PASSIVE+ANL_INDICATIVE+ANL_PAST} ja_{CO-ORDINATED_CONJUNCTION} punnittiin_{CO-ORDINATED_VERB+THINK}. *Selvässä asiassa*_{LOCATION+NOTION} ei_{NEGATIVE-AUXILIARY+ANL_NEGATION+ANL_THIRD+(ANL_SINGULAR)} ole_{ADJACENT_AUXILIARY} enempää_{QUANTITY+MUCH} pohdimista_{ANL_INFINITIVE4}. *Artikkeli*_{AGENT+COMMUNICATION, ANL_OVERT} pohti_{(ANL_ACTIVE)+ANL_INDICATIVE+ANL_THIRD+(ANL_SINGULAR)} kysymystä_{PATIENT+COMMUNICATION/(NOTION)}, onko_(PATIENT+INDIRECT_QUESTION) --.

Table 2.11. Approximate English translation of the lexical entry for *pohtia* in *Perussanakirja* (PS) and its linguistic analysis; default lexical entry forms (i.e., sentence-initial FIRST INFINITIVES) in parentheses.

pohtia^{61*F}

think about something_{PATIENT+NOTION?} thoroughly_{MANNER+THOROUGH}, [evaluating different possibilities]_{MANNER+THOROUGH}, consider, [think], deem, think, reason, weigh, riddle.

*Ponder*_(INFINITIVE1) a riddle_{PATIENT+NOTION/COMMUNICATION} a problem_{PATIENT+NOTION}.

*Ponder*_(INFINITIVE1) the question_{PATIENT+COMMUNICATION} from every angle_{MANNER+THOROUGH}.

*Ponder*_(INFINITIVE1) means_{PATIENT+ACTIVITY/(NOTION)} to help_{PURPOSE/REASON+ACTIVITY} in a matter.

Table 2.12. Approximate English translation of the lexical entry for *pohtia* in *Nykysuomen sanakirja* (NS) and its linguistic analysis; default lexical entry forms (i.e., sentence initial FIRST INFINITIVES) as well as default features (i.e., ACTIVE voice and SINGULAR number) in parentheses.

<p>pohtia^{17*} (verb)</p> <p>1. (=winnow) </p> <p><i>The grain</i>_{PATIENT+SUBSTANCE} <i>was threshed</i>_{ANL_PASSIVE+ANL_INDICATIVE+ANL_PAST} [<i>with a thresher</i>]_{INSTRUMENT+ARTIFACT} --</p> <p>(usually) 2. consider, [think], deem, think, reason, weigh, riddle </p> <p>Ponder_(INFINITIVE1) <i>some matter</i>_{PATIENT+NOTION}, <i>situation</i>_{PATIENT+STATE}.</p> <p>Ponder_(INFINITIVE1) <i>the means</i>_{PATIENT+ACTIVITY/(NOTION)} <i>to help</i>_{PURPOSE/REASON+ACTIVITY} <i>in some matter</i>.</p> <p><i>The question</i>_{PATIENT+NOTION/COMMUNICATION} was pondered_{ANL_PASSIVE+ANL_INDICATIVE+ANL_PAST} <i>and</i>_{COORDINATED_CONJUNCTION} <i>weighed</i>_{CO-ORDINATED_VERB+THINK}.</p> <p><i>In a clear matter</i>_{LOCATION+NOTION} not_{NEGATIVE-AUXILIARY+ANL_NEGATION+ANL_THIRD+(ANL_SINGULAR)}</p> <p><i>is</i>_{ADJACENT_AUXILIARY} [<i>there</i>] <i>more</i>_{QUANTITY+MUCH} pondering_{ANL_INFINITIVE4}.</p> <p><i>The article</i>_{AGENT+COMMUNICATION, ANL_OVERT}</p> <p>pondered_{(ANL_ACTIVE)+ANL_INDICATIVE+ANL_THIRD+(ANL_SINGULAR)} <i>the question</i>_{PATIENT+COMMUNICATION/(NOTION)}, <i>whether</i>_(PATIENT+INDIRECT_QUESTION) --.</p>

With respect to the syntactic and semantic argument context for *pohtia* manifested in only these very small sets of examples, we can already start to see some emergent characteristics, consisting prominently of different types of PATIENT arguments representing ABSTRACT NOTIONS (including STATES), ACTIVITIES, and forms of COMMUNICATION. We can further scrutinize the overall occurrence and frequencies of these patterns by combining them all together with regard to each dictionary source, presented in Table 2.13 for *pohtia*. In this aggregate representation, I have decided to exclude contextual information from the multiword definitions as well as default forms and features, namely, the FIRST INFINITIVE when used as the solitary head of an example fragment without an auxiliary FINITE verb, as well as the ACTIVE voice, which always applies for any FINITE form with a person/number feature, and SINGULAR number which predominates among the ACTIVE FINITE forms of the studied THINK lexemes, since these do not in my opinion convey any essential additional characteristic differentiating information.²⁹ Furthermore, this consideration of these particular features as default characteristics will be motivated in the selection of contextual variables for inclusion in the multivariate analysis which will follow later in Sections 3.4.2 and 5.1. We can now again see that NS contains a larger range and more occurrences of contextual information in comparison to PS. As far as morphological features are concerned, the two sources have in common only the default FIRST INFINITIVE form, but among syntactic arguments and their semantic classifications, the aforementioned three abstract semantic types as PATIENTS as well as PURPOSE (or alternatively interpreted as REASON) as arguments have persisted from NS to PS as characteristic of *pohtia*. It is worth noting that the only new syntactic argument type for *pohtia*, which is present in PS but not in NS, is the THOROUGH

²⁹ One could very well ask why the INDICATIVE mood is not also considered as a default feature here, since this could be argued to be the case on the basis of the examples for *pohtia*. However, it will turn out with the other three THINK lexemes that another mood, namely, the IMPERATIVE, also has several occurrences, so I have consequently decided to keep both of these two moods as part of the analysis. With respect to the SINGULAR vs. PLURAL opposition in number, there is only one single PLURAL FINITE form among the scrutinized examples, which in my opinion is not enough to warrant the marking of all SINGULAR forms.

subtype of MANNER; this may have arisen with the abstractization of the meaning of the word.

Table 2.13. Aggregated linguistic analysis of the lexical entry example sentences for *pohtia* in both *Perussanakirja* (PS) and *Nykysuomen sanakirja* (NS); default lexical entry forms (i.e., sentence-initial FIRST INFINITIVES) as well as default features (i.e., ACTIVE voice) in parentheses; common features in **boldface**.

Contextual features/ <i>pohtia</i>	PS	NS
NEGATION	0	+
INDICATIVE	0	++
PAST	0	+
(ACTIVE)	0	(+)
PASSIVE	0	+
THIRD	0	++
OVERT	0	+
(SINGULAR)	0	(++)
(INFINITIVE1)	(+++)	(++)
INFINITIVE4	0	+
AGENT +COMMUNICATION	0	+
PATIENT + NOTION +STATE +ACTIVITY +COMMUNICATION +INDIRECT QUESTION	+ 0 + + 0	+ + + ++ (+)
MANNER +THOROUGH	+	0
QUANTITY +MUCH	0	+
LOCATION +NOTION	0 0	+ 0
PURPOSE/REASON (+ACTIVITY)	+	+
VERB-CHAIN +NEGATIVE_AUXILIARY +ADJACENT_AUXILIARY	0 0	+ +
CO-ORDINATED CONJUNCTION	0	+
CO-ORDINATED_VERB +THINK	0	+

This linguistic analysis process of the usage examples can now be replicated for the entire set of the studied THINK lexemes (plus *tuumia/tuumata*, which was ruled out solely on the basis of its relatively lower frequency), which are presented in Appendix G. Together, these yield the overall results presented in Table 2.14. It would be tempting to apply the battery of statistical analyses to be presented later in this dissertation in Section 3 to this dictionary content data, but the observed frequencies (as they stand) are far too low to produce even remotely reliable results, so we must content ourselves with a general qualitative description. Overall, the two dictionaries contained exemplars of 26 morphological or related features pertaining to the inflected form or morpho-syntactic role of the studied THINK lexemes themselves, 55 couplings of a syntactic argument and their semantic subclassifications, and 4 unclassified syntactic argument types. Of these, 9 had an occurrence with all of the four studied THINK lexemes in at least one of the two sources, and two in both, namely,

abstract NOTIONS and ACTIVITIES as PATIENTS. The other contextual features common to all four THINK lexemes in at least one of the sources are the INDICATIVE mood, the PASSIVE voice, the PAST tense, and the THIRD person, and the OVERT subject/AGENT among the morphological features, the THOROUGH subtype of MANNER, and ADJACENT AUXILIARY verbs as part of the verb chain. All of these can be broadly understood as prototypical of neutral dictionary entries.

As could be expected, the number of senses lexeme-wise correlates with the range of the exemplified possible contexts, so that there are 68 contextual feature associations for *ajatella*, 47 for *mieltiä*, 39 for *harkita*, and 19 for *pohtia*. Feature-wise, the contents of the PS corresponds for the most part with NS, in that NS has as many or (often substantially) more exemplars of the co-occurrence of some contextual feature and a particular lexeme. However, there are 16 cases where there are more co-occurrences of a feature and some lexeme in PS than in NS, and 13 of these were not observable at all in NS. Nevertheless, for all but one of these there is only a singular exemplar in PS. Thus, these comparisons of these two dictionaries indicate differences broadly reminiscent of the findings made by Atkins and Levin (1995: 90-96) concerning the treatment of English near-synonymous *shake* verbs in three English language-learner dictionaries. The most abundantly exemplified syntactic argument is the PATIENT (PS:33; NS:72), followed by the VERB-CHAIN (PS:23; NS:40), the AGENT (PS:16; NS:51), and MANNER (PS:10; NS:19), with the other argument types clearly trailing behind. What is most interesting is that morphological features of the verb (including the verb-chain it may be part of), practically ignored in AR, are exemplified considerably more than are individual syntactic arguments, with joint frequencies of 115 in PS and as many as 295 in NS.

Table 2.14 Lexeme-wise aggregates of the occurrences of the selected contextual features in the linguistic analyses of the example sentences for the four studied THINK lexemes in both *Perussanakirja* (PS) and *Nykysuomen sanakirja* (NS), with the first value indicating the frequency of occurrences in PS and the second value that in NS; default lexical entry forms (i.e., sentence-initial FIRST INFINITIVES) as well as default features (i.e., ACTIVE voice and SINGULAR number) are not considered; features with occurrences in conjunction with all four THINK lexemes underlined; features with occurrences with all but one of the four THINK lexemes ~~struck through~~; features with occurrences with only one lexeme in either source in **boldface**; features with more occurrences per one or more lexemes in PS than NS in *italics*. In addition, the occurrences of contextual features are presented for the *tuumia/tuumata*, but these figures are not included in the just-mentioned assessments, and features present only in the usage examples of *tuumia/tuumata* but none of the studied four THINK lexemes are in (parentheses).

Contextual features/Lexemes	<i>ajatella</i>	<i>mieltiä</i>	<i>pohtia</i>	<i>harkita</i>	(<i>tuumia/tuumata</i>)
MORPHOLOGY (115..295)					
+NEGATION	1..4	0..2	0..1	-	-
+INDICATIVE	<u>9..22</u>	<u>2..15</u>	<u>0..2</u>	<u>5..8</u>	(4..17)
+IMPERATIVE	2..4	0..2	-	-	-
+PRESENT	3..8	2..1	-	3..2	(1..3)
+PAST	<u>3..11</u>	<u>0..10</u>	<u>0..1</u>	<u>1..2</u>	(3..13)
+PASSIVE	<u>5..3</u>	<u>0..1</u>	<u>0..1</u>	<u>1..5</u>	(0..2)
+FIRST	3..7	0..3	-	3..3	(0..3)
+SECOND	3..4	1..2	-	-	(1..4)
+THIRD	<u>3..19</u>	<u>3..12</u>	<u>0..2</u>	<u>5..4</u>	(3..10)
+PLURAL	0..1	-	-	-	(1..0)

+ <u>OVERT</u>	<u>1..19</u>	<u>0..6</u>	<u>0..1</u>	<u>1..4</u>	(2..10)
+ <i>COVERT</i>	7..10	4..9	-	5..2	(2..7)
+INFINITIVE1	2..8	1..0	-	1..2	(1..2)
+INFINITIVE2	1..2	-	-	1..2	-
+INFINITIVE3	1..2	1..2	-	-	(0..1)
+ <i>INFINITIVE4</i>	<i>1..0</i>	-	-	1..0	(0..1)
+PARTICIPLE1	3..4	0..3	-	1..6	(0..1)
+PARTICIPLE2	4..6	1..4	-	3..9	(0..3)
+ ESSIVE	-	-	-	0..1	-
+TRANSLATIVE	1..1	-	-	2..0	-
+ INESSIVE	1..1	-	-	-	-
+ ELATIVE	0..1	-	-	-	(0..1)
+ ILLATIVE	-	0..1	-	-	-
+ABESSIVE	1..1	1..1	-	-	-
+INSTRUCTIVE	1..2	-	-	1..2	-
+ <u>CLAUSE_EQUIVALENT</u>	5..9	1..4	-	4..8	(0..2)
AGENT (16..51)					
+ INDIVIDUAL	8..26	4..15	-	3..4	(4..17)
+GROUP	0..1	-	-	1..2	-
+ BODY	0..1	-	-	-	-
+ ARTIFACT	-	0..1	-	-	-
+ COMMUNICATION	-	-	0..1	-	-
PATIENT (33..72)					
+ INDIVIDUAL	2..3	1..1	-	1..1	(1..2)
+ <i>INDIVIDUAL</i>	2..5	-	-	-	-
+ <i>FAUNA</i>	<i>1..0</i>	-	-	-	-
+ ARTIFACT	0..1	-	-	-	(0..1)
+LOCATION	1..3	-	-	0..1	-
+ <u>NOTION</u>	<u>6..4</u>	<u>1..3</u>	<u>1..1</u>	<u>4..7</u>	(1..3)
+ <i>STATE</i>	-	-	0..1	<i>1..0</i>	-
+ATTRIBUTE	0..1	0..1	-	0..1	(0..1)
+ <i>TIME</i>	<i>1..0</i>	-	-	-	-
+ <u>ACTIVITY</u>	<u>1..5</u>	<u>1..2</u>	<u>1..1</u>	<u>1..5</u>	-
+COMMUNICATION	-	1..3	1..2	0..2	-
+ COGNITION	-	0..2	-	-	-
+INFINITIVE1	1..3	0..4	-	1..2	(1..7)
+INDIRECT_QUESTION	1..3	1..1	-	0..1	-
+ DIRECT_QUOTE	-	1..1	-	-	(1..4)
+ <i>että 'that' clause</i>	<i>2..1</i>	<i>0..1</i>	-	-	-
SOURCE (1..2)					
+ INDIVIDUAL	0..1	-	-	-	-
+ NOTION	1..1	-	-	-	(1..1)
GOAL (4..12)					
+INDIVIDUAL	0..2	0..1	-	-	(0..1)
+ <i>NOTION</i>	<i>1..0</i>	-	-	-	-
+ATTRIBUTE	0..1	0..1	-	2..5	-
+ LOCATION	1..2	-	-	-	-
MANNER (10..19)					
+ GENERIC	-	-	-	1..1	(0..2)
+ <i>POSITIVE (CLARITY)</i>	<i>2..1</i>	-	-	0..1	-
+ NOTION/ATTRIBUTE	1..1	-	-	-	-
+ <u>THOROUGH</u>	<u>0..2</u>	<u>1..1</u>	<u>1..0</u>	<u>0..5</u>	-
+ <i>CONCUR</i>	<i>1..0</i>	-	-	-	-
+ DIFFER	1..1	-	-	-	(0..1)
+ALONE	0..1	0..1	-	-	(0..1)

(+TOGETHER)	-	-	-	-	-
+ FRAME	1..1	-	-	-	-
+ LIKENESS	0..1	-	-	-	-
+ ATTITUDE	-	-	-	0..2	-
+ SOUND	1..0	-	-	-	-
(+TIME)	-	-	-	-	(0..1)
(COMITATIVE)	-	-	-	-	(0..1)
QUANTITY (1..5)					
+ <i>MUCH</i>	-	<i>1..0</i>	0..1	0..2	-
+LITTLE	0..1	0..1	-	-	(0..1)
LOCATION (0..3)	-	0..2	-	-	-
+ NOTION	-	-	0..1	-	-
+EVENT	-	-	-	-	(0..1)
TMP (1..5)					
+ INDEFINITE	0..2	1..2	-	0..1	(0..2)
DURATION (0..4)					
+OPEN	0..1	0..1	-	-	(0..2)
+ LONG	0..2	-	-	-	(0..1)
+SHORT	-	-	-	-	(0..2)
PURPOSE/REASON (1..3)	-	0..1	1..1	0..1	-
(META [Clause-Adverbial])	-	-	-	-	(0..2)
VERB-CHAIN (23..40)					
+ NEGATED_AUXILIARY	1..3	0..2	0..1	-	-
+ ADJACENT_AUXILIARY	<u>8..12</u>	<u>3..3</u>	<u>0..1</u>	<u>3..4</u>	(1..5)
+ <i>COMPLEMENT</i>	<i>1..0</i>	-	-	0..1	(0..1)
+PROPOSSIBILITY	2..3	-	-	0..1	-
+ IMPOSSIBILITY	0..1	-	-	-	-
+ PRONECESSITY	1..2	1..1	-	1..3	(0..1)
+ TEMPORAL	-	0..1	-	-	-
+ <i>CAUSE</i>	<i>1..0</i>	-	-	-	-
+ ACCIDENTAL	1..1	-	-	-	-
CO-ORDINATING CONJUNCTION (0..2)	-	0..1	0..1	-	(0..2)
CO-ORDINATED_VERB (0..5)					
+ THINK	0..1	0..1	0..1	-	-
+ COGNITION	-	0..1	-	-	-
+ VERBAL	0..1	-	-	-	-
(+ACTION)	-	-	-	-	(0..2)

In any case, the individual frequencies of the contextual features among the examples for such an extremely limited data in terms of its size are less important than their occurrences or nonoccurrences in conjunction with each studied THINK lexeme. Thus, the key observation at this stage is that the examples do indicate clear differences in the usage of the studied THINK lexemes: 17 (20.0%) of the altogether 85 possible contextual features did not exhibit a co-occurrence with *ajatella*; the corresponding non-co-occurrence figures are 38 (44.7%) for *miettiinä*, 65 (76.5%) for *pohtia*, and 46 (54.1%) for *harkita*. Furthermore, 35 (41.2.8%) of all the contextual features had a co-occurrence with only one of the studied lexemes in either dictionary (presented by each lexeme in Table 2.15 below), and 10 (11.8%) of these singular preferences were consistent in both sources. These latter features cluster around *ajatella*, associated with the INESSIVE case, human INDIVIDUALS as PATIENT, abstract NOTION as SOURCE, physical LOCATION as GOAL, NOTION/ATTRIBUTE, DIFFER, and FRAME as MANNER, as well as an ACCIDENTAL verb chain for *ajatella*, while among the three other lexemes

only *mieltä* is consistently associated with a DIRECT QUOTE as PATIENT and *harkita* with the GENERIC type of MANNER. In contrast, there were 19 (22.4%) features which had occurrences with all but one of the studied THINK lexemes, constituting a type of negative evidence by way of absence (presented by each lexeme in Table 2.15 below); among these, 6 (7.1%) features had such (non-)occurrence patterns in both dictionaries. These latter absences of a feature in comparison to the three other lexemes focused all for *pohtia*, being PRESENT tense, COVERT subjects, PAST (SECOND) PARTICIPLE, CLAUSE-EQUIVALENT usage, human INDIVIDUALS as AGENT, and positive NECESSITY (i.e., obligation) in the verb-chain. Interestingly, there are no contextual features for which *mieltä* would be the only one of the studied THINK lexemes without an occurrence among the dictionary usage examples.

Table 2.15. Lexeme-wise sole occurrences and sole absences, in contrast to the three other THINK lexemes at a time, of contextual features among the usage examples in PS and NS; occurrences in both sources in **boldface**, occurrences only in PS underlined, occurrences only in NS in *italics*.

Lexeme/ Feature	Sole occurrences	Sole absences
ajatella	<i>PLURAL</i> , INESSIVE , <i>ELATIVE</i> , <i>AGENT+BODY</i> , PATIENT+INDIVIDUAL , <u><i>PATIENT+FAUNA</i></u> , <u><i>PATIENT+ARTIFACT</i></u> , <u><i>PATIENT+TIME</i></u> , <i>SOURCE+INDIVIDUAL</i> , SOURCE+NOTION , <u><i>GOAL+NOTION</i></u> , GOAL+LOCATION , MANNER+NOTION/ATTRIBUTE , <u><i>MANNER+CONCUR</i></u> , MANNER+DIFFER , MANNER+FRAME , <i>MANNER+LIKENESS</i> , <u><i>MANNER+SOUND</i></u> , <i>DURATION+LONG</i> , <i>VERB-CHAIN+IMPOSSIBILITY</i> , <u><i>VERB-CHAIN+CAUSE</i></u> , VERB-CHAIN+ACCIDENTAL , <i>CO-ORDINATED VERB+VERBAL</i>	PATIENT+COMMUNICATION, PURPOSE(/REASON)
miettiä	<i>ILLATIVE</i> , <i>AGENT+ARTIFACT</i> , <i>PATIENT+COGNITION</i> , PATIENT+DIRECT_QUOTE , <i>LOCATION(+GENERIC)</i> , <i>VERB-CHAIN+TEMPORAL</i> , <i>CO-ORDINATED VERB+COGNITION</i>	-
pohtia	<i>AGENT+COMMUNICATION</i> , <i>LOCATION+NOTION</i>	NEGATION, PRESENT , FIRST, COVERT , INFINITIVE1, PARTICIPLE1, PARTICIPLE2 , CLAUSE_EQUIVALENT , AGENT+INDIVIDUAL , <i>PATIENT+ATTRIBUTE</i> , PATIENT+INFINITIVE1, PATIENT+INDIRECT_QUESTION, GOAL+ATTRIBUTE, TMP+INDEFINITE, VERB-CHAIN+PRONECESSITY
harkita	<i>ESSIVE</i> , MANNER+GENERIC , <i>MANNER+ATTITUDE</i>	VERB-CHAIN+NEGATED_AUXILIARY, CO-ORDINATED_VERB+THINK

What will be my interest vis-à-vis these contextual features in this dissertation is the extent to which their co-occurrences or absences with the studied THINK lexemes in the two dictionaries will correspond with their actual usage in the extensive corpus data. Furthermore, I aim to order these features in terms of their relative importance for each lexeme with the help of the same data.

2.3.3 The etymological origins of the selected THINK lexemes

I will now close this overview of the present, existing descriptions of the studied THINK lexemes with a look back into their past, turning to what is currently known of their etymology. Excerpts translated into English of the latest explanations for the origins of these lexemes according to *Suomen sanojen alkuperä* (SSA) by Itkonen, Kulonen et al. (1992-2000) are presented in full in Appendix H. Of the four lexemes in question, one is a complex derivative of an old Finnic root with a hunting-related meaning, while two are abstractions of originally rural/agricultural verbs with concrete activities as their referents, and only one has apparently been loaned with its present cognitive meaning largely intact.

The most common of the set *ajatella*, is believed to be a frequentative further derived form of the FACTIVE (CAUSATIVE) derivation *ajattaa* of the verb *ajaa* ‘drive/chase’. It is conceived of having been originally understood as the figurative “chasing” and pursuit of the object of thought, still used in this meaning in, for example, *ajan takaa* ‘I am driving/chasing after/from behind’, which can still also be seen to mean *koetan palauttaa tai saada mieleeni* ‘I am trying to recall or get [something] back into my mind’, or, alternatively, ‘my [ultimate] intention is ...’. In turn, the root *ajaa* may possibly be an Indo-European loan. In its current meaning ‘think’ *ajatella* is quite opaque to the average native speaker of Finnish with limited knowledge of etymology – such as myself, prior to this study – with respect to its morphological and semantic constitution, thus conforming with Fortescue’s (2001: 30) conclusion concerning languages in general in this respect. Nevertheless, *ajatella* can easily be seen as a derivation using still fully productive elements (i.e., the causative *-ttA-* followed by the frequentative *-ele*, e.g., Karlsson 1983: 201, see also 2008) in Finnish and a current root, that is, *ajaa*, when one is pointed in the right direction.

For its part, *harkita* still also means (or has recently meant) in many Finnish dialects the quite concrete activity of *harata*, *naarata jotakin veden pohjasta* ‘trawl/drag something from the bottom of a body of water’ in addition to the more abstract, cognitive meaning. It can be derived (with the productive morpheme *-tA-*, e.g., Karlsson 1983: 201, see also 2008) from the noun *harkki*, meaning a variety of mostly countryside-related referents, for example, ‘twig/branch harrow, dragnet; fork-headed spade for lifting potatoes; fork-headed hay pole; fork/branch; a type of device for weaving nets’ and many more, but to my understanding this meaning has diminished with the urbanization of Finland. Likewise, *pohtia* was not long ago seen primarily as a parallel form of the quite concrete farming activity *pohtaa* ‘winnow’, specifically to separate the wheat from the chaff, as is still exemplified in NS. Though to many native Finns of the older generations with roots in the countryside *pohtia* may still appear as a relatively transparent metaphorical extension of meaning, similar to *punnita* as both ‘weigh’ and ‘consider’, to myself and others of the younger generation with a purely urban background the underlying more concrete denotation is no longer commonly accessible. This verb is considered either to have a descriptive origin, or alternatively to be a loan into Early Proto-Finnic from Pre-Germanic. Finally, *mieltiä* is the only one in the quartet believed to have been borrowed in more or less in its current meaning. With respect to its original source, two explanations have been suggested. The one considered more probable traces *mieltiä* to Slavic, corresponding to Russian *smétit* ‘guess, assume, notice, grasp/understand’, while a

secondary association is assumed via Estonian *mõtelda* ‘think < *mõõta* ‘measure’ in the Germanic root **mēt-*, corresponding to the modern Swedish *mäta* ‘measure’.

Thus, three of the most common Finnish THINK lexemes have their origins in rural life, in hunting (i.e., *ajatella*), fishing (i.e., *harkita*), and farming (i.e., *pohtia*), though these associations and related meanings have become increasingly synchronically opaque or peripheral for most native speakers of modern Finnish (for a sketch of their modern usage, see Länsimäki 2007; for a popularized overview of this general rural characteristic of modern Finnish words and expressions, see Repo 2003). Therefore, Fortescue’s (2001: 30-31) assessment that the most basic verbs of THINKing would generally stem in languages mostly from more visible/perceivable [mental] states and activities, such as speaking/pronouncing, observing, or wishing/intending, would not appear to hold in the case of Finnish, unless one extends the possible scope of origins further back to include the actual *physical* activities from which these more abstract senses are derived. Furthermore, Fortescue’s (2001: 29, Example 5) listing (2.2 below) of the most common metaphorical expressions underlying verbs of THINKing seems in this light incomplete, as it lacks THINKing as *searching/seeking/chasing/hunting* (i.e., *ajatella* and *harkita*) and THINKing as sifting and separating apart (with considerable toil, i.e., *pohtia*), evident in these Finnish THINK lexemes.³⁰ Firstly, Fortescue (2001: 28) sees ‘finding’ rather as a case of polysemous extension of THINK lexemes (2.3 below) than as a possible origin from which their present COGNITIVE meaning might have been metaphorically abstracted. Secondly, the uncontested metaphorical origins of *pohtia* and *harkita*, evident also in the English ‘barnyard’ terms ‘brood’ and ‘ruminant’, are in Fortescue’s (2001: 30-31) view secondary, evaluative and culture-specific in nature, an assessment which would not appear to hold in the case of Finnish due to the high relative frequency and semantic generality of these two THINK lexemes.

(2.2) Polysemies of THINKing

- a. thinking = *believing* ~ *being true/truthful* | *saying/pronouncing* (~ *hearing*)
- b. thinking = *considering/judging* ~ *being true/truthful* |
saying/pronouncing ~ *finding* (~ *hearing*)
- c. thinking = *unspecified/general mental activity* (~ *hearing*)
- (d). thinking = *intending*

(2.3) Metaphorical expressions underlying THINKing

- a. thinking < *weighing*
- b. thinking < *observing*
- c. thinking < *wanting*
- d. thinking < *calculating*
- (e). thinking < *worrying*

³⁰ Of course, one could conceive of the searching/seeking aspect of *ajatella* in the abstract sense to denote INTENTION, and thus fall under (2.3c).

2.4 The compilation of the research corpus and its general description

2.4.1 General criteria

In contrast to the early days of corpus linguistics, the size and variety of electronic corpora available to linguistic research has grown tremendously over the last few decades, and even more so with the World Wide Web and other electronic media (for a concise summary of this development, see Kilgarriff and Grefenstette 2003: 334-335, 337-340). With for instance the 180 million word Finnish Text Collection (FTC 2001) as only one of the existing resources for already a few years, a researcher of Finnish does have some choice and does not have to resort to a pure convenience sample. However, because transcribed and annotated spoken language resources are still very limited for Finnish, the range of choices considered in this dissertation is restricted to written corpora. Within this mode of language use, the research corpus used in this study was compiled in accordance with several guiding principles. In general, these selectional criteria should be *external*, that is, social and contextual, and thus essentially not based solely on the linguistic content (Clear 1992; Sinclair 2005).

Firstly, as prior linguistic research has indicated that individual speakers or writers do have individual preferences, which can have at least some influence on the results (e.g., Bresnan et al. 2007), I decided to use corpora in which the writer or speaker of each text fragment is consistently identifiable. Furthermore, it was desirable that the number of writers or speakers who took part in producing the corpus at any particular time were substantial – in the order of several hundreds with respect to the studied THINK lexemes – in order to be able to take into account and dilute the influence of overtly idiosyncratic individuals or individual instances of idiosyncratic usage. Secondly, it follows from the first principle that it would be desirable to have (many) more than one observed usage of the studied THINK lexemes from as many as possible of the identified writers in order to be able to study individual consistency (or inconsistency) as well as idiolectal preferences concerning the usage of the studied lexemes. Although I have opted for a larger number of authors, the extent of exemplars per each writer (on the average) is restricted in that the linguistic analysis is validated and supplemented manually. As the individual outputs of a large number of writers or speakers can be considered more independent of each other than the more lengthy output sequences of a few or only one person, this also fits with general statistical sampling theory (Woods et al. 1986: 104-105). Thirdly, it was my intuition that this would entail temporally coherent, contiguous stretches of corpora, instead of unconnected random samples (referred to as collections of *citations* rather than as a “proper” corpus by Clear 1992), even more so as this would also allow for the later study of intratextual cohesion and repetitiveness between separate texts produced by the same author around approximately the same period of time.

Fourthly, I wanted to study and describe contemporary Finnish usage which was at the same time conformant with the general norms and conventions of standard written Finnish (i.e., *kirjakieli* ‘book Finnish’ with a focus on word-by-word orthography rather than punctuational correctness), in order to allow for its automatic parsing, but also language which is nevertheless produced in and for the moment at hand, if not genuinely fully spontaneous, and is thus not heavily edited, or otherwise repeatedly considered, reviewed, and polished. The former criterion would on the one hand rule out text messages and other recent types of electronic telegraph-style

communications, where orthographical and other rules are bent due to the limitations of available space (see, e.g., Kotilainen 2007a, 2007b, forthcoming, for studies of Finnish using material from a variety of such “new” electronic registers and genres, including web pages, weblogs and chat forums, in addition to the newsgroup discussion to be included in this study, as well as Kukko 2003 for text messages via mobile phones), but on the other hand, also fiction and non-fiction book-length literature such as novels or scientific text books. Nevertheless, my aim was to study and describe what is considered by non-linguist native speakers as “good and presentable” Finnish usage. In practice, the latter criterion would mean a preference for a large number of shorter texts, produced and published within a day or so, over a small number of, or singular, longer texts from each individual writer, which may have been worked on for longer periods of time.

However, since prior research has also shown both variance as well as cohesion and repetition effects within individual texts written by individual authors (e.g., Hoey 1991; Thompson 1998), as the fifth principle, all the individual texts from the corpora to be selected would be included in their entirety, this also being a practice recommended by Sinclair (2005). Together with the third criterion of temporal contiguity, this would also allow for not only the later study of interrelationships between texts produced by different individuals concerning the same topic around the same time, but also the intratextual relationships of fragments by different individuals within the same text, for example, direct citations within newspaper articles and (possibly recursive) quotations within Internet newsgroup discussion postings. Sixthly, I wanted to use general sources which are inherently heterogeneous and diverse with respect to the topics and subjects they cover, even though I would not use or need all the available material within this study. Nevertheless, I would rather focus on and cover comprehensively only a small number of such sources, which would furthermore be clearly distinct from each other (i.e., again a form of scientific *triangulation*), rather than attempt to canvass a wide range of different sources and registers, genres and text types. Thus, I do not attempt to compile a generally *balanced* research corpus, which in the view of many would in any case be a difficult if not impossible task (Atkins et al. 1992).

Finally, the exact size of the selected corpus material would be determined quite simply by how long a contiguous sequence of basic subunits from the selected sources would exhibit a sufficient number of all four selected THINK lexemes and their distinct contexts. This, in turn, was influenced by the requirements of the statistical methods to be discussed later in Section 3, and, in practice, meant several hundred occurrences for the least frequent lexeme, and several thousands of occurrences for the selected lexemes altogether. With respect to the contextual features, the quantitative sufficiency of the selected corpus samples could be assessed through the extent that adding more subunits of data would substantially introduce instances of new, previously unobserved features, stopping at a point when the growth of possible variation could clearly be judged to have reached a plateau (similar to the “freezing point” referred to in Hakulinen et al. 1980). One should note here that some of these criteria were set primarily rather to allow for the re-use of the material to be analyzed here in later textual or discourse-oriented analyses, than due to obligatory requirements arising solely from this current study. Two corpus sources which fit the above criteria were 1) newspapers and 2) Internet newsgroup discussions, among others.

As representatives of these two sources for use in this dissertation, I selected two months (January–February 1995) of written text from Helsingin Sanomat (1995), Finland’s major daily newspaper, and six months (October 2002– April 2003) of written discussion in the SFNET (2002-2003) Internet discussion forum, primarily concerning (personal) relationships (`sfnet.keskustelu.ihmissuhteet`) and politics (`sfnet.keskustelu.politiikka`). The general characteristics of the text types incorporated in these two sources are discussed in depth in Appendix I. Furthermore, Appendix I will also contain a detailed description of the various stages in the selection and processing of samples from these two sources for inclusion in the actual research corpus, as well as the structural make-up and demographic and other characteristics of this corpus.

2.4.2 Main characteristics of the final selected research corpus

The contents of the final research corpus resultant after the considerations and processing presented in Appendix I are described in Table 2.16. In the final newspaper subcorpus, there were 1323 articles containing an occurrence of one or more of the studied THINK lexemes, divided into 1007 articles with exactly one single occurrence and 316 articles with two or more, the maximum being 11 within one individual article. With respect to the identity of writers, 296 journalists or otherwise identifiable authors used the studied lexemes at least once, of which 87 authors exactly once and 219 twice or more (the maximum being 27 for an identifiable author in 17 articles), while a slight majority of 230 did not use the studied lexemes at all in the selected newspaper subcorpus. In the final Internet newsgroup subcorpus, there were 1318 postings which contained at least one occurrence of the studied THINK lexemes, divided into 1085 postings with only a single occurrence and 233 postings with two or more instances, with a maximum of 9 occurrences within a single posting. Among the individual identifiable contributors, 251 used the studied lexemes at least once, of which 83 exactly once and 168 twice or more, with a maximum of 146 over 93 postings (for contributor #721), while a clear majority of 922 did not use any of the studied lexemes even a single time.³¹

³¹ One must remember here that the presented figures pertaining to the overall number of contributors in both sources, and thus the proportions of those having used or not used the THINK lexemes, can only be considered approximate, depending on which identity codes are included in the calculations, that is, do we consider as contributors only those who can be exactly identified, having both a first and a last name, or those for whom only their gender or native-language status can be established, or simply all distinct recorded author designations (newspaper article author codes and newsgroup posting e-mail addresses). In the ensuing statistical analyses, for practical purposes all distinct author codes which I have not been able to combine with sufficient certainty are understood to refer to an individual author, whether pertaining to an exactly identifiable individual or not.

Table 2.16. Figures describing the final research corpus and its two component subparts.

Statistics/subcorpus	HS	SFNET
THINK lexemes	1750	1654
Words (including punctuation)	4011064	1400020
Words (excluding punctuation)	3304512	1174693
Individual texts	16107	18729
Individual texts with identifiable authors	10569	(18729)
Individual texts containing THINK lexemes	1323	1318
Individual texts with THINK lexemes and an identifiable author	1049	1318
Individual identifiable authors using THINK lexemes	296	251
Usage of THINK lexemes with an identifiable author	1392	-

The frequencies of the individual THINK lexemes in each subcorpus are presented numerically in Table 2.17 and graphically in Figure 2.5. As can be seen, the selected lexeme quartet is again clearly more frequent than the rest in both partitions of the research corpus. However, in comparison to the rankings in the FTC (2001), *ajatella* is now overall the most frequent of the group, and this is also the case in both subcorpora, though in the newspaper material *pohtia* is a close second, reminiscent of the ranking order in FTC. Furthermore, the frequency range for the studied four THINK lexemes is narrower in the newspaper text than in the Internet newsgroup discussions, where the frequency differences are somewhat more pronounced. Finally, it is interesting to note the occurrence, though yet quite infrequent, of the compound forms *toisinajatella* ‘differ disagree/think differently’, *samoinajatella* ‘agree/concur/think similarly’ and *pitkäänmieltä* ‘think/ponder long’ in the research corpus, which might be indications of commencing lexicalization of some of the contextual associations (examples of which can be found among the prominent argument-specific lexemes in Table P.1 in Appendix P). Moreover, looking at Figure 2.5, we can also see that the distribution of the THINK lexemes in the research corpus overall is not exactly Zipfian, with the four most frequent lexemes having observed frequencies clearly above the ideal Zipfian distribution, while the number of observations of the less frequent lexemes in contrast fall below this ideal. Furthermore, the ratio (0.748) of the most common THINK lexeme, here unequivocally *ajatella* in contrast to the case with FTC, against all the other THINK lexemes together now exceeds the rough equality that Manin (submitted) has hypothesized.

Table 2.17. Frequencies of the selected four THINK lexemes as well as the other, less frequent ones in both subcorpora; lexemes in (parentheses) are novel compound constructions outside the original lexeme set presented in Section 2.1.4.

Lexeme/frequency	Newspaper subcorpus (HS)	Internet newsgroup discussion subcorpus (SFNET)	Research corpus altogether
ajatella	570	922	1492
mieltä	355	457	712
pohtia	556	157	713
harkita	269	118	387
punnita	45	13	58
tuumia	41	7	28
mietiskellä	17	7	24
aprikoida	12	2	14
hautoa	11	6	17
järkeillä	6	9	15
tuumata	8	3	11
filosofoida	4	7	11
funtsia	1	1	2
funtsata	1	1	2
(toisin#ajatella)	1	1	2
(samoin#ajatella)	0	1	1
(pitkään#mieltä)	0	1	1

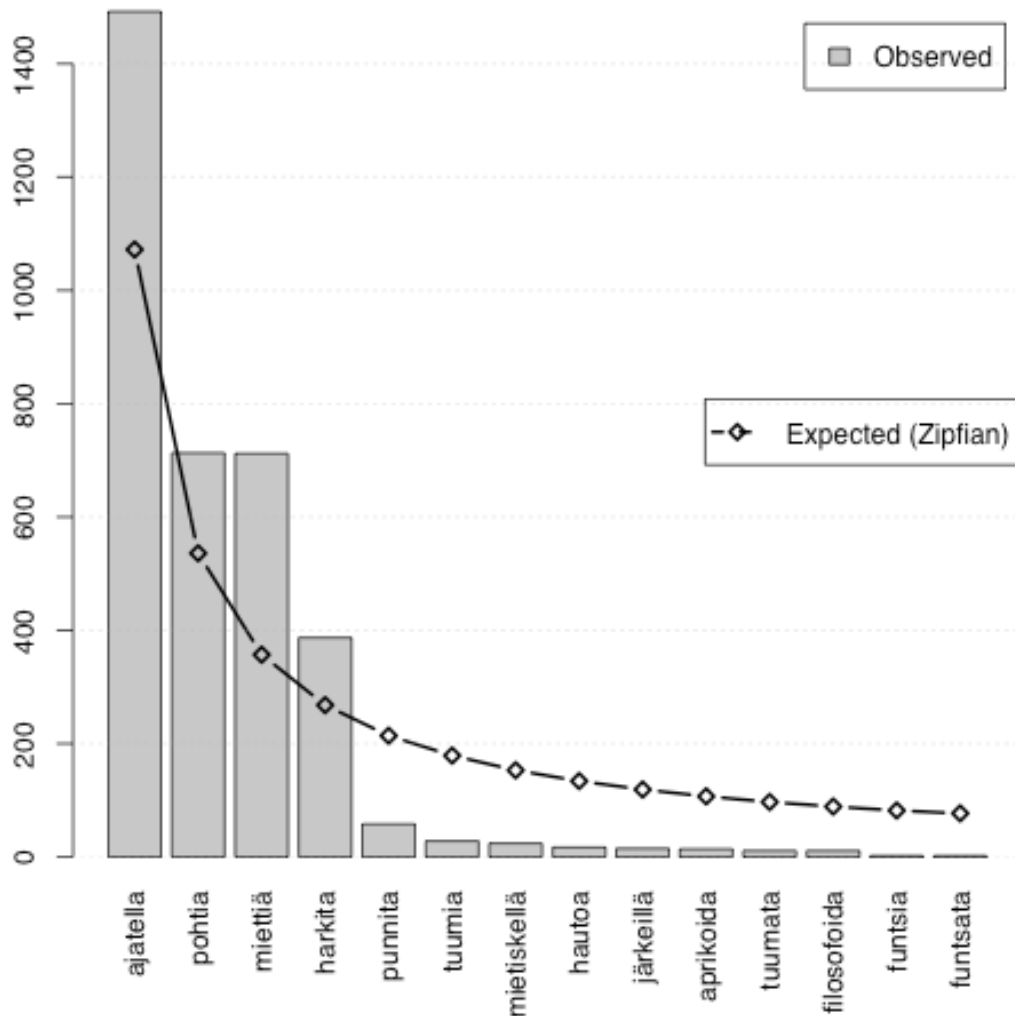


Figure 2.5. Frequencies of the THINK lexemes in Table 2.17, contrasted with an ideal Zipfian distribution based on their joint frequencies.

2.4.3 Coverage of contextual features in the research corpus

At this stage we can assess the sufficiency of the selected quantities of research corpus with respect to the studied linguistic features. I will firstly study the accumulation of new morphological features and their clusters, that is, entire inflected forms, for the studied THINK lexemes, as well as the frequency of the most infrequent lexeme in the synonym set, namely, *harkita*, both in the three distinct portions of the research corpus, i.e., the newspaper subcorpus and the two newsgroups, and the research corpus as a whole (Figures 2.6-2.9). If we look at the occurrences of new morphological features in these Figures, we can see that their number reaches a plateau of approximately 40 distinct features (approximately four-fifths of the total of 52 features theoretically possible in verbal word forms) or so in all the subcorpora, and subsequently also in the entire corpus, exhibiting just such a curvilinear distribution as one would expect for type frequency (Biber 1993: 185). In the newspaper material as well as the relationships newsgroup this happens by the end of

the first quarter of the subcorpora in question, and a similar trend appears to also apply to the number of morphological features with at least two occurrences (which Sinclair 2005 considers as a minimum evidence recurrence to be considered as an independent linguistic event), that in the newspaper corpus ceases to grow at the beginning of the second quarter, while the relationships newsgroup takes somewhat longer to reach this stage, but nevertheless clearly before the end of the second quarter.

In contrast, one has to go well into the second half of the politics newsgroup before the increase in the number of new morphological features flattens out, and for the proportion of these features with at least two occurrences to reach the same level requires almost the entire content of this particular newsgroup, clearly longer than is the case in the overall similar-sized relationships newsgroup. We must remember, however, that the two newsgroup portions are in overall size approximately one-sixth each in comparison to the newspaper material, though both sources contain roughly as many THINK lexemes in absolute terms. So, if we map on top of each other the growth rates of morphological features with at least two occurrences in the three portions (Figure 2.10), we can see that the number of new features grows faster in both of the two newsgroups than in the newspaper material, with the latter taking much longer in terms of running text to reach the overall maximum plateau. Consequently, the Internet newsgroup subcorpus can be considered more “rich” in THINK lexemes.

As for the accumulation of new inflected forms, their number continues to grow steadily throughout all these three corpus portions, though the growth rate of forms with at least two occurrences clearly slows down after the initial surge by the end of first half of each subcorpus. However, the slope steepens slightly at roughly the juncture in the sequential make-up of the research corpus where the newspaper subcorpus ends and the newsgroup subcorpus begins, which then returns to slower but still steady growth. This point of discontinuity in the growth-rate curve can be considered indicative of some level of difference among the two subcorpora. Regarding the occurrences of *harkita*, this particular lexeme seems to be quite evenly dispersed in the three corpus portions, exhibiting roughly a linear growth rate which is to be expected for token frequencies (Biber 1993: 185), though there appear to be some dry zones especially in the politics newsgroup, at the very beginning and at approximately the two-thirds milestone.

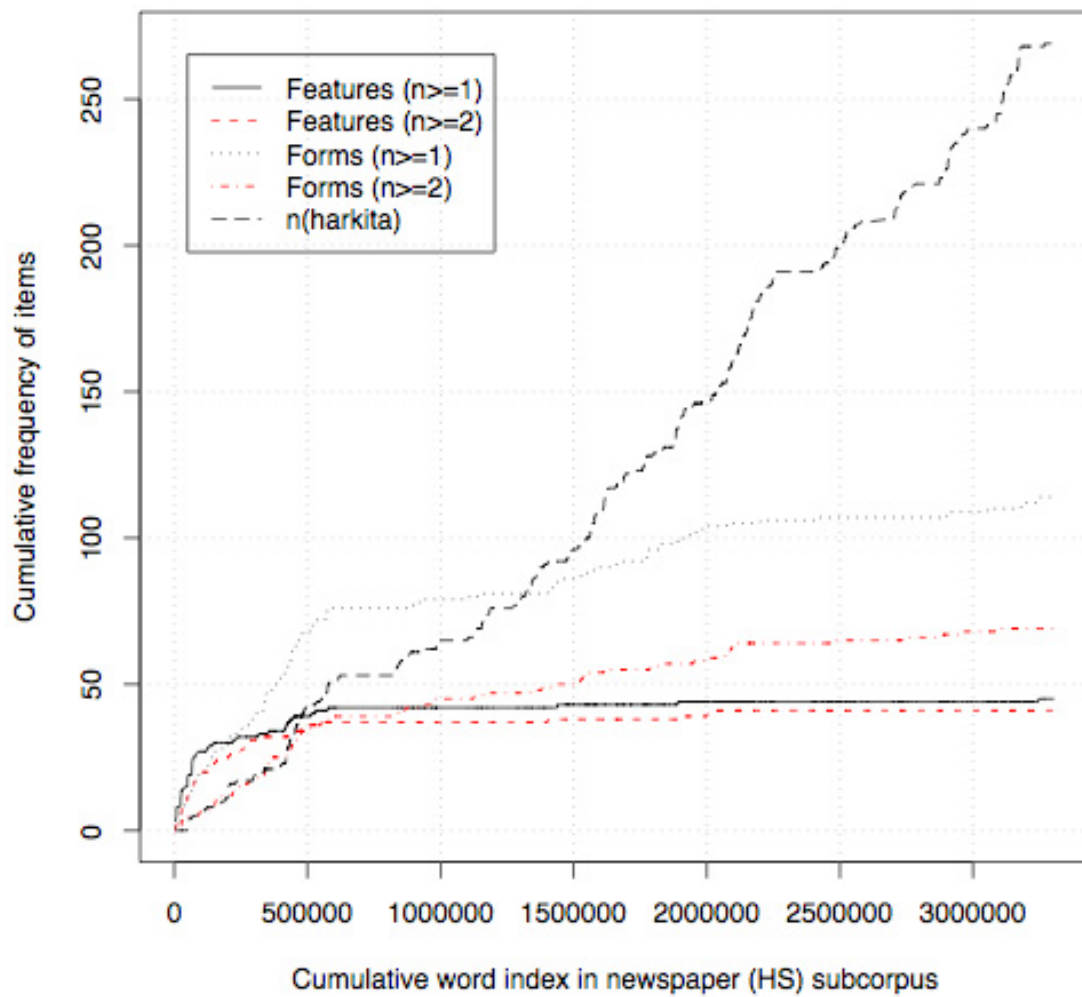


Figure 2.6. Growth rates of the individual morphological features and their clusters as distinct inflected forms, as well as the occurrences of *harkita* in the newspaper (HS) subcorpus; with a distinction between at least one and at least two observations of each scrutinized type.

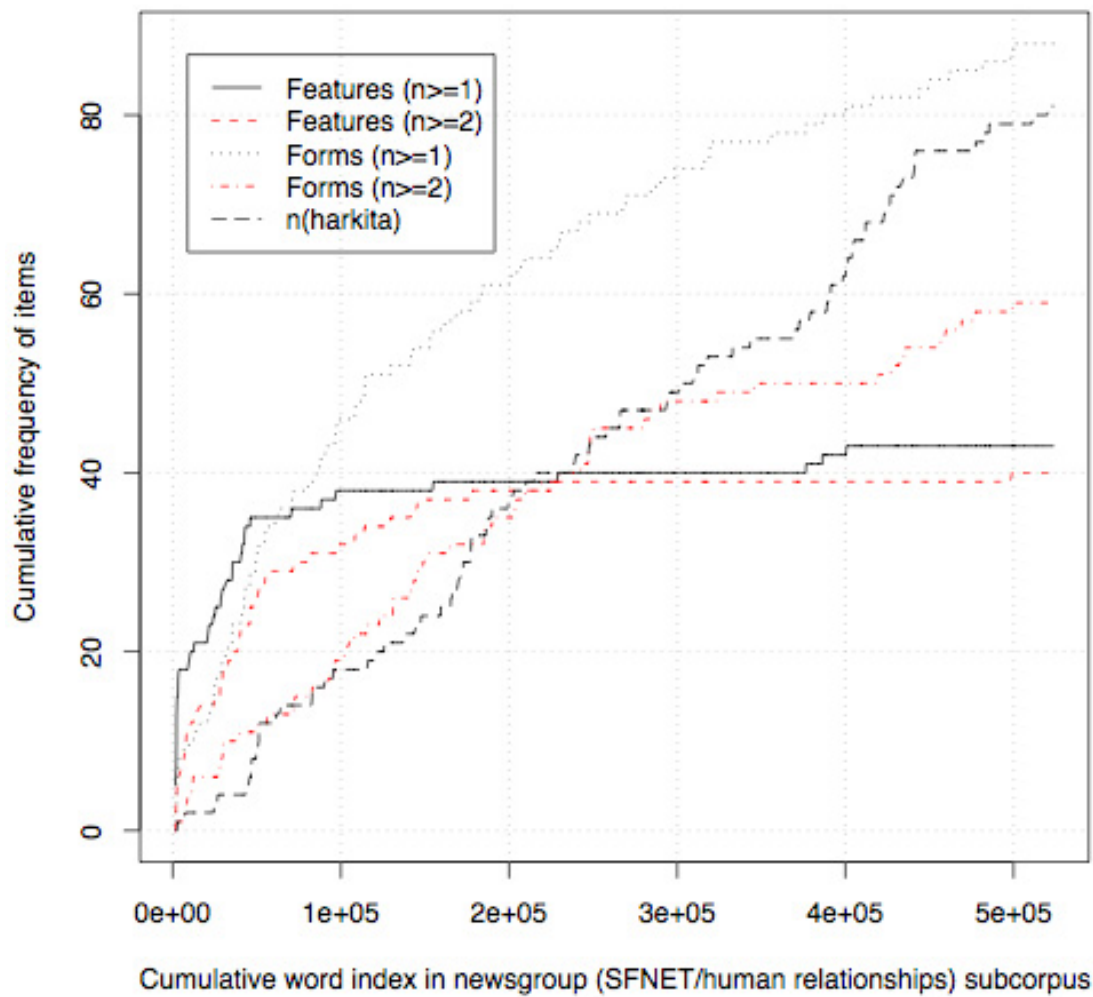


Figure 2.7. Growth rates of the individual morphological features and their clusters as distinct inflected forms, as well as the occurrences of *harkita*, in the `relationships` newsgroup portion of the SFNET subcorpus; with a distinction between at least one and at least two observations of each scrutinized type.

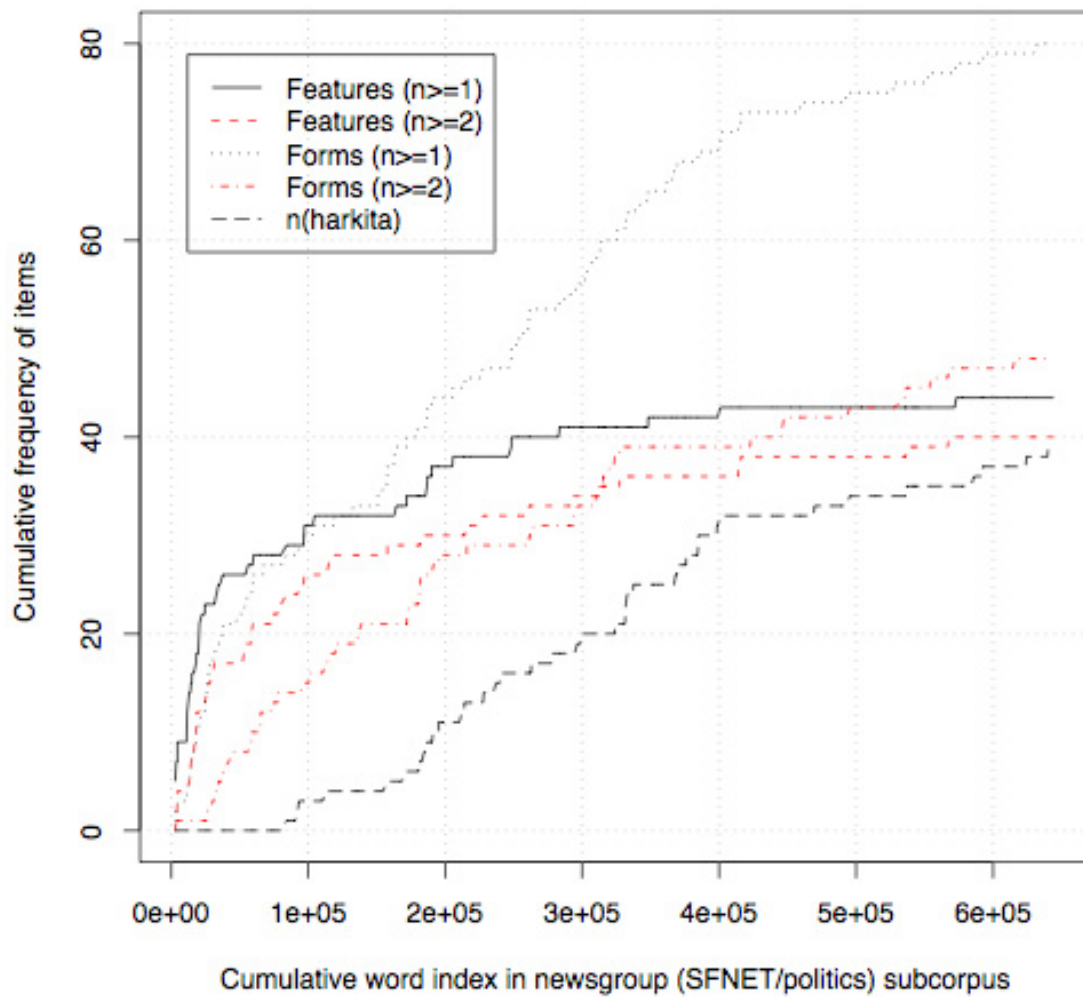


Figure 2.8. Growth rates of the individual morphological features and their clusters as distinct inflected forms, as well as the occurrences of *harkita*, in the `politics` news group portion of the SFNET subcorpus; with a distinction between at least one and at least two observations of each scrutinized type.

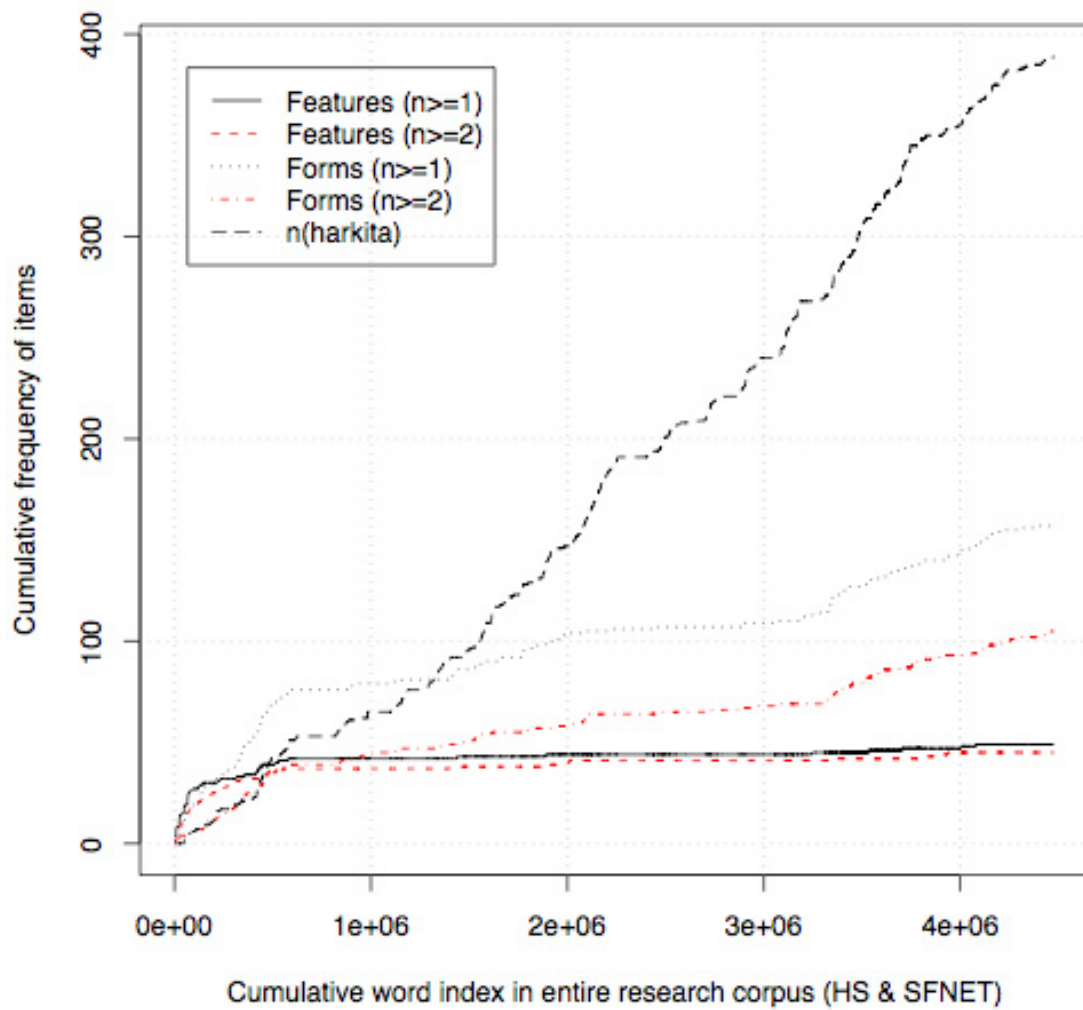


Figure 2.9. Growth rates of the individual morphological features and their clusters as distinct inflected forms, as well as the occurrences of *harkita* in the entire research corpus; with a distinction between at least one and at least two observations of each scrutinized type.

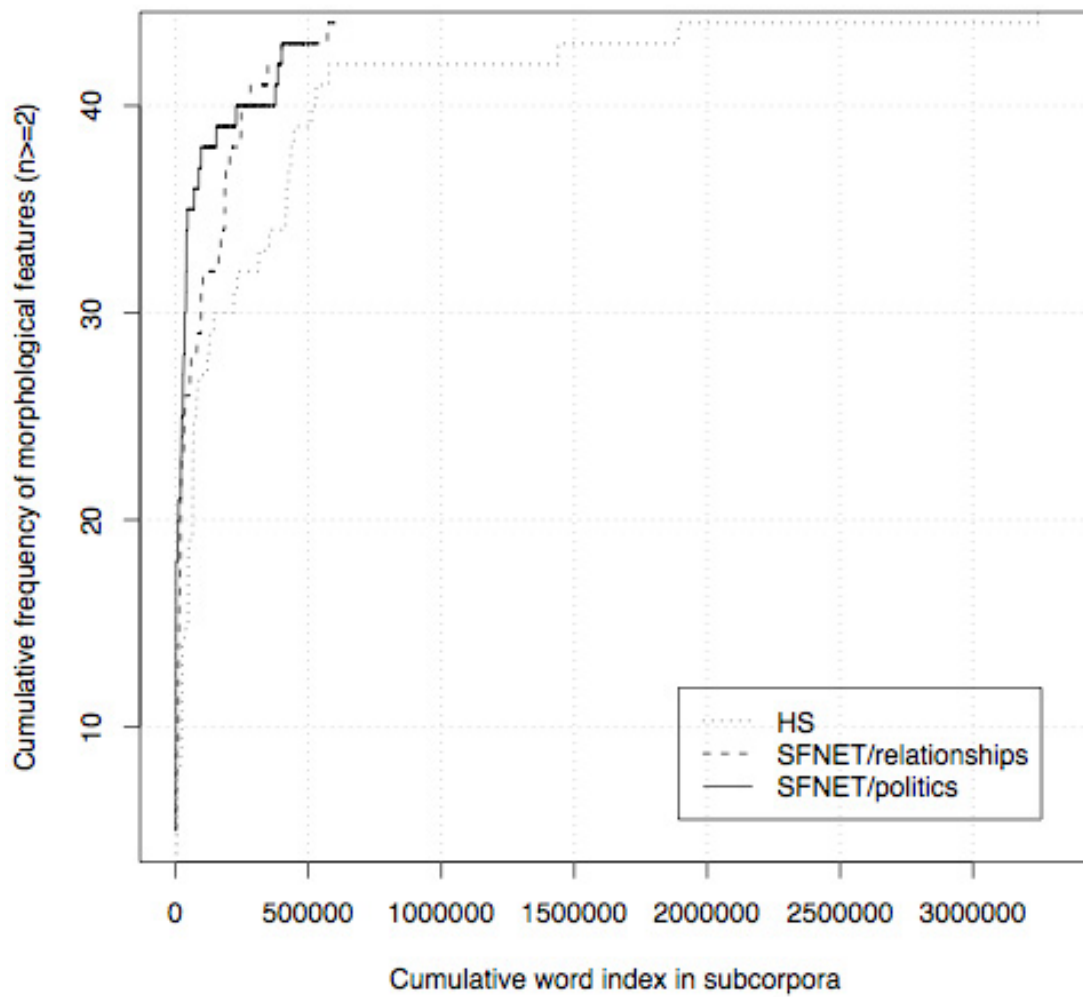


Figure 2.10. Growth rates of the number of individual morphological features with at least two occurrences in each of the three distinct subcorpora.

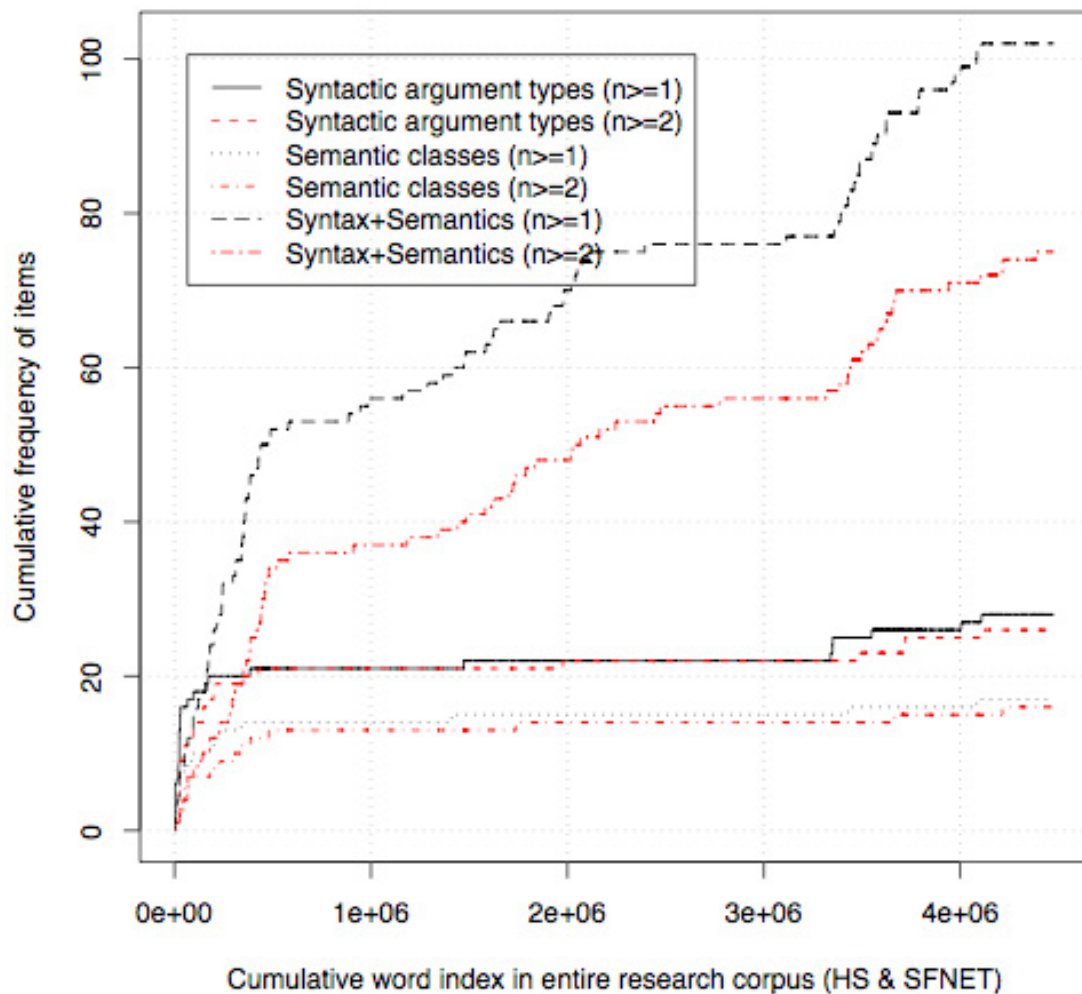


Figure 2.11. Growth rates of the syntactic argument types and their semantic classification types (restricted to nominal arguments, according to WordNet) as well as the combinations of syntactic arguments and semantic classes in the entire research corpus; with a distinction between at least one and at least two observations of each scrutinized type.

Secondly, I will observe the rate at which the number of distinct syntactic argument types and the semantic classifications of nominal arguments (following the 25 unique beginners for nouns in WordNet) grow in the entire research corpus (Figure 2.11). It takes only a small fraction of the entire research corpus for us to have observed at least one occurrence of the syntactic argument types existent in its entire content, and reaching at least two occurrences for this set follows soon thereafter. Interestingly, there appears to be a small surge towards the last quarter of the research corpus, which is roughly the point where the Internet newsgroup subcorpus starts. Keeping the sequential structure of the research corpus in mind, this implies that the newsgroups would appear to contain a few new syntactic structures in comparison to the newspaper material. Regarding the growth rate of the observations of distinct semantic classes, this is somewhat more gradual, but also plateaus in terms of both the first and second occurrences of each type at approximately the same point in the

corpus as is the case with the syntactic arguments. Likewise, there seems to be a small notch upwards at the point where the newsgroup material begins.

The growth of the combinations of syntactic arguments and nominal semantic classes appears analogous to that of the inflected forms above, with an initial, fast surge followed by a less steep but continually ascending slope. At roughly half-way through the entire corpus, which is towards the end of the newspaper portion, the growth rate practically flattens out. This is followed by a noticeable second surge where one can assume the newsgroup material begins, which then eases again down to a gentler slope, thus, in practice, repeating the prior development stages in the newspaper portion of the research corpus. This would suggest some structural differences between the two sources included in the research corpus, but it remains to be seen whether these differences will also be reflected in the frequency counts as statistically significant.

Third, one could then evaluate whether the growth of the occurrences of individual features is stable with respect to their proportions among the studied THINK lexemes. In practice, this is a worthwhile exercise, as long as the overall frequencies of the features in question are sufficiently high in the entire research corpus, at the minimum several tens and preferably at least one hundred, so that there is enough data to exhibit visually observable trends. In principle, one could assess all the sufficiently frequent features, but for reasons of space I decided to scrutinize as examples the FIRST PERSON SINGULAR among the node-specific morphological features (with 248 occurrences altogether in the research corpus), and human GROUPS as AGENTS among the combinations of syntactic and semantic features (with 256 occurrences), since these two have been the object of previous studies (Arppe 2002; Arppe and Järviö 2007b).

As we can see in Figures 2.12-2.13, the overall trends of both features in relation to the studied THINK lexemes appear to change at the boundary of the two subcorpora, although within each subcorpus the growth rates appear quite stable. In fact, the two features behave quite differently in the two subcorpora, so that in comparison to the newspaper text the FIRST PERSON SINGULAR merely increases its growth rate in the Internet newsgroup material (though it is hard to tell exactly for *pohtia* and *harkita* which are proportionately quite insignificant in comparison to *ajatella* and *mieltä*; after an initial spurt in the *relationships* newsgroup *harkita* does not appear to occur in conjunction with FIRST PERSON SINGULAR at all in the *politics* newsgroup.) Overall, this would be in line with the higher density of THINK lexemes in the newsgroup subportion, which, nevertheless, at 1.3 instances per 10,000 words (i.e., $150 \cdot 10000 / 1174693$) is cross-linguistically very low, in comparison to 35/10000 reported for English, 2.6 for Swedish, and 9 for Dutch (Goddard 2003: 132), even though it combines the occurrences of all four of the THINK lexemes.³²

With respect to GROUPS as AGENT, their lexeme-specific growth rates show changes at the subcorpus boundary, so that the strongest growth of occurrences that this feature exhibited with *pohtia* in the newspaper text turns in the newsgroup portion into an

³² If we consider only the occurrences of *ajatella*, the proportion is reduced to as low as 0.8/10000 (i.e., $97 \cdot 10000 / 1174693$). However, if we include all the cases for the four THINK lexemes in which the verb-chain as a whole exhibits the FIRST PERSON SINGULAR feature (if not the lexemes themselves), the resultant proportion is similar to that reported for Swedish at 2.6 (i.e., $309 \cdot 10000 / 1174693$).

effective standstill. In contrast, *ajatella*, which had exhibited the lowest proportion of occurrences with GROUP AGENTS in the newspaper text, picks up speed and reaches by the end of the newsgroup portion the same level in absolute terms as *mieltä*, which had kept its growth rate relatively stable throughout, similar to *harkita*. This, too, could be interpreted as a clear indicator of linguistic differences between the two subcorpora. In general, the above observations are a reminder for us that there can be, or rather, most probably will be variation between genres, something which one cannot ignore in our analysis (cf. Biber 1998).

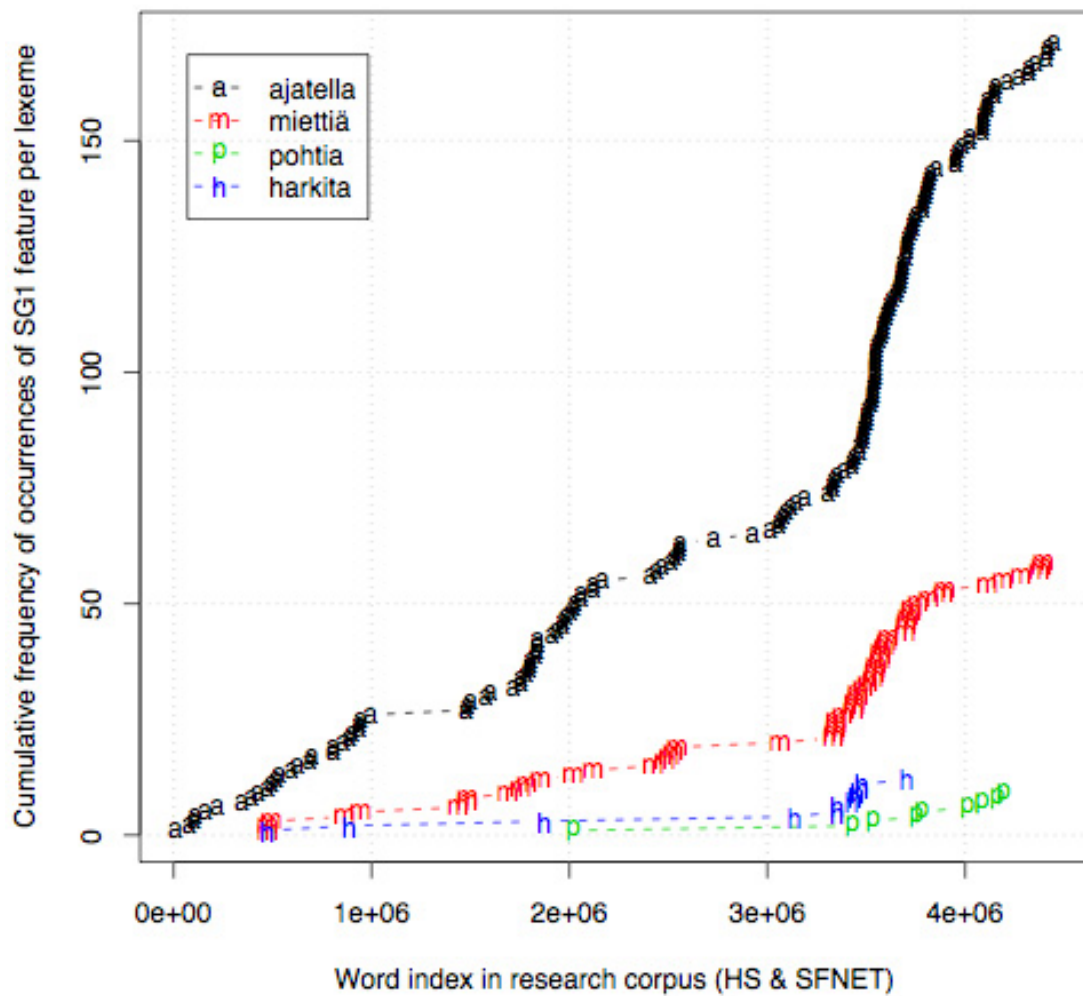


Figure 2.12. Growth rate of the FIRST PERSON SINGULAR feature among the studied THINK lexemes in the research corpus.

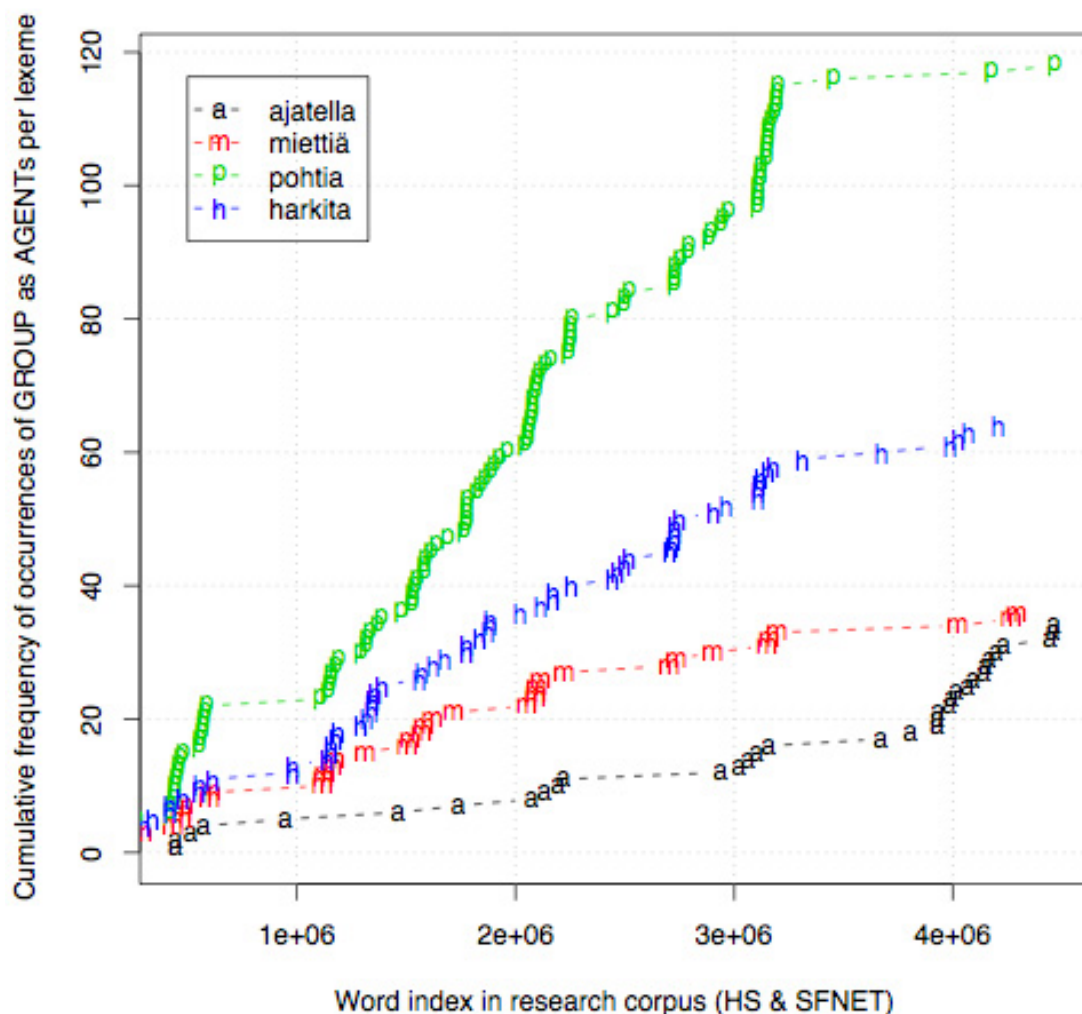


Figure 2.13. Growth rate of human GROUPs as AGENTS of the studied THINK lexemes in the research corpus.

In conclusion, on the basis of these scrutinies of the contextual feature content of the research corpus as well as the individual subcorpora, I consider the selected samples as sufficiently large to cover at least the most frequent and typical feature contexts of the studied THINK lexemes with respect to the selected sources/media. This is based on a visual examination of diminishing increments in line with what Biber (1983: 190) recommends with respect to curvilinear growth-rates particular to feature types, such as are studied in this dissertation. Furthermore, as long as we remember to pay attention to the distinctions between the two subcorpora and the genres they represent, the overall proportions of individual features appear to exhibit relatively stable proportions among the studied THINK lexemes. Finally, one should note that since the research corpus and its subsections are not genuinely random extracted samples, the statistical calculations suggested by Biber (1993: 183-195) concerning sufficient sample size are neither applicable nor relevant here, even more so as I have no prior data (based on some other representative sample) concerning the necessary initial estimates of the expected variances for the parameters of interest.

2.4.4 Representativeness, replicability, reliability, and validity

In empirical linguistic research based on corpora, it has become difficult to avoid the question of the *representativeness* of the material being used with respect to the general phenomenon under study here, namely, language, being in this dissertation specifically contemporary Finnish. Representativeness in the context of statistical inference (e.g., Woods et al. 1986: 48-58, 77-94; Howell 1999: 5-8, 21-22) inherently entails that we can define and demarcate a general overall *population* of entities or events relevant and specific to the particular scientific field and object of research, from which we can then take a *sample* (N.B. concerning in the statistical context *measures* of *values* of the characteristics that interested us), a kind of crude but truthful snapshot. Such sampling is motivated only when it is impractical or impossible to grasp, record, and study the entire population as it is defined; if we can with relative ease cover all the possible and relevant entities in the population, there is no reason to use a sample in its stead. The sample can be considerably smaller than the entire population; as far as we compile the sample either *randomly* or according to criteria and proportions that *accurately* reflect the entire population (i.e., *stratified* sampling), the sample will represent the properties of the entire population within a range of accuracy determined by statistical sampling theory. In such a situation, one can on the basis of the sample alone make *generalizations* concerning the entire population, that is, the sample is then representative of the population and the phenomena that it incorporates and which are measured. Otherwise, if the aforementioned requirements are not met, the characteristics of the sample may in the worst case be limited to reflect with certainty only the sample itself, in other words, the results are not generalizable.

In the end, linguists such as myself using corpora wish to make general statements about the entire linguistic system of whose productive output the contents of corpora are. The difficulty in this is, firstly, that language as the multimodal and multidimensional physical, biological, psychological, and social human phenomenon that it is in the broad sense does not constitute a well-defined and uniform population, as Kilgarriff and Grefenstette (2003: 340-341) very convincingly illustrate (see also, e.g., Atkins et al. 1992: 4-5; Leech 2007: 134-136). Secondly, as a natural consequence of this multifacetedness of both language as an individual and a collective human phenomenon and of its interpersonally observable incarnations, that is, texts and utterances in whatever physical form they may take, there are no obviously clear-cut, concrete *units* of language which one could use as a basis for compiling a sample. Leech (2007: 138) proposes as such a sampling unit the abstraction *atomic communicative event* (ACE), which is a trio constituted by the communicated linguistic fragment, its initiator and each of its receivers individually, but it is clear that just the linguistic content of such an ACE will take many forms and representations.

In general, it appears to me that corpus-oriented linguists have for the most part been fixated with considering only recorded or recordable corpora as relevant linguistic evidence (exemplified by Sampson 2005³³), and the problems of principle concerning

³³ Not only does Sampson (2005: 17-18, 28) consider elicitation as no different from native-language introspection by the linguist him/herself, but he also views experimentation as data concerning linguistic *feelings* or judgments, for which there is “no *a priori* reason to assume that such data must be

representativeness, and the discrepancy between what one wants and purports to study and what is actually observed and described, are restricted to and revolve around corpora and corpora alone. That is, either corpora that are currently at one's disposal or novel corpora (or extensions of existing ones) which one can in practice get hold of or create within a reasonable time. Then, as a way out of this theoretical snag one may as the first option take a pragmatic stance and make do with the corpora that one has, but at the same openly acknowledge the associated shortcomings and limitations concerning the interpretation of the results (e.g., Clear 1992: 21-22, 31; Atkins et al. 1992: 4-5; Manning and Schütze 1999: 120; Kilgarriff and Grefenstette 2003: 334; see also Stubbs 1996: 232).³⁴ As a second alternative, one may resort to extralinguistic (sociocultural) criteria and professional (possibly collective) judgement to select an individual *exemplary corpus* (Bungarten 1979: 42-43, cited in Leech 2007: 137). More comprehensively, one may aim to enumerate all the possible (situationally and functionally defined) categories of language use such as genres and registers and estimate their relative proportions, and then compile a corresponding sample for each type, producing a *balanced corpus* (Biber 1993).

The problem with the latter approach of stratified sampling is that the selection of categories for inclusion into the corpus and especially the estimation of their sampling proportions is, when accommodating their cultural importance rather than actual occurrence (if such could be at all reliably estimated, as Atkins et al. [1992: 6] note), normative at best and utterly subjective at worst, and does not result in and correspond to what is generally considered as statistically representative sampling (Váradi 2001: 590-592). Furthermore, Biber's explicit purpose is to cover the "full range of linguistic variation existing in a language ..." instead of "summary statistics for the entire language [represented in the corpus]", which generalizations are in his mind "typically not of interest in linguistics" (Biber 1993: 181). I am not convinced that this is universally the case; in contrast, my interest is similar to that in the COBUILD project, namely, "the central and typical uses of the language" (attributed to Patrick Hanks in Clear et al. 1996: 304), which we will, in the results in Sections 4-5, of this dissertation see to contain ample variation in itself in the case of the selected THINK lexemes. Thirdly and finally, one may expect that continuing to increase the size and diversity of corpora will by itself start to alleviate their alleged lack of representativeness (e.g., Clear 1992: 30; Kilgarriff and Grefenstette 2003: 336; Leech 2007: 138).

Personally, I take the position that there is no point in this attempt to reduce all the distinct aspects of language use into a one and the same, all-encompassing corpus, or a fixed set of such corpora, for that matter. In terms of corpus-based study, one should rather, in my view, tackle the complex multifacetedness of language piece by piece, by picking individual distinct types of linguistic usage (i.e., corpora) and covering these comprehensively one at a time, developing and testing hypotheses about the general underlying linguistic system gradually along the way. More generally speaking, in order to really understand language as the multimodal phenomenon that

a reliable guide to the properties of a speaker's language, ...", despite arguments to the contrary (see, e.g., Arppe and Järviö 2007a).

³⁴ An extreme example of paying lip-service to this issue is exhibited by Manning and Schütze (1999: 120): "In summary, there is no easy way of determining whether a corpus is representative, but it is an important issue to keep in mind when doing Statistical NLP work", leaving the matter at that, and preceded by "... and one should simply use all the text available."

it is, one should cover the various ways and processes through which language is conceived of, produced, communicated, received, and understood, and make the most of the different types of linguistic evidence and methods that are presently available, in addition to corpora, for example, traditional elicitation, experimentation, and the like (e.g., Arppe and Järvi­kivi 2007a, 2007b). This, in fact, is the pluralistic view of linguistic research that Chafe (1992) has already argued for quite some time. Thus, I advocate a shift in the focus: from worrying about whether one particular study is fully representative of (all) language (use) to whether the results can be repeatedly replicated in heterogeneous settings, be this via divergent corpora or through entirely different research methods. Incidentally, this approach to validating results is what statisticians such as Rosnow and Rosenthal (1989) and Moran (2003) currently recommend, rather than putting all efforts and energy into increasing sample sizes in individual studies.

Consequently, I do not claim my research corpus to be representative of all Finnish in the strictest statistical sense, but of the genres to which the two subcorpora belong, and by extension, perhaps also contemporary written Finnish in general. Nevertheless, I believe the two subcorpora of newspapers and Internet newsgroup discussion to both be exemplary corpora in the spirit of Bungarten (1979: 42-43, as cited in Leech 2007: 137). I justify this view on the external grounds that newspapers are as a textual genre and form of communication considered by sociologists and communication researchers (Groth 1960: 125; Pietilä 1997: 43) as a central societal glue in a contemporary (Western) society, which Finland certainly is. Furthermore, Helsingin Sanomat, the newspaper in question, can in particular be characterized as culturally important on the objective basis, in accordance with Leech (2007: 139), due to its extensive readership in Finland noted in Appendix I.2. In turn, the Internet newsgroup discussion subcorpus can also be accorded a special position as it consists first and foremost of interpersonal human conversation, which Chafe (1992: 88-89) argues as the most basic kind of human linguistic interaction (though strictly speaking he is referring to speaking). Likewise, Biber (1993: 181) does concede that conversation probably accounts for the great majority of all actual linguistic usage, estimating its proportion as high as 90%, as does in approximate terms also Clear (1992: 24-26) with respect to language production.

A natural continuation of the corpus-based results to be presented in this study is to test and to try to replicate them with various types of experiments, along the lines as was undertaken with respect to a subset of the selected THINK lexemes and contextual feature variables in an earlier study (Arppe and Järvi­kivi 2007b). This runs in contrast to simply increasing the number of different text types which are covered or increasing the sizes of samples scrutinized, as Leech (2007: 138) would appear to suggest as the next step forward. In the aforementioned prior study, when journalists in Jyväskylä writing newspaper articles, and engineering students from all around Finland sweating in an exam in Espoo, as well as regulars and occasional patrons in a Helsinki pub developing or recovering from a hangover, both groups while taking a few minutes to respond to an experimental questionnaire, all produced convergent linguistic evidence, we considered the overall result as sufficiently and convincingly revealing of the particular phenomenon in Finnish. With respect to the considerably more complex setting scrutinized in this study, I would look to a similar multimethodologically rigorous validation of the ensuing results.

2.5 Preparation of the corpus for statistical analysis

The corpus data used in this dissertation was annotated and manipulated in several stages using an assortment of UNIX shell scripts written by me, roughly sketched in 2.4 below. These scripts – plus the original corpus data and the subsequent linguistic analyses – are all available in the microcorpus *amph*, located under the auspices of CSC – Center of Scientific Computing, Finland, to be found at <URL: <http://www.csc.fi/english/research/software/amph>>. The linguistic content of the original two subcorpora was first automatically analyzed using the FI-FDG parser at stage (2.4a), while leaving the extralinguistic structure and mark-up intact. The resultant morphological and syntactic analyses of the studied THINK lexemes and their contexts in the corpora were verified and supplemented by hand at stage (2.4b), at which time also the semantic classification of nominal arguments was undertaken. After this, the ensuing morphological, syntactic, semantic and phrase-structural analyses of the occurrences of the studied lexemes and their context, within the selected portions of the original data, including both the identified syntactic arguments and simple linear context of five words both to the left and right of the node, as well as extra-linguistic data present in the non-linguistic structure of the original corpora, were extracted in several stages (2.4c-g), which were then transformed into a text-format data table suitable for the *R* statistical programming and computing environment at stage (2.4h), with the occurrence of some feature in a context marked as TRUE and its absence as FALSE.

At this point, verb-chain-specific analytical tags were added, and in the spirit of Arppe (2002) all possible permutations of up to three features (of any kind) were generated for each extracted lexeme, whether in the context of the studied lexemes or one of the studied THINK lexemes themselves. A large majority of such feature combinations would turn out to be singular occurrences so that they will become redundant by any statistical test or cut-off frequency, but this full scale application of combinatorics allows for the possibility of the most common and possibly statistically most significant combinations or underlying sub-combinations of features, that is, abstractions of patterns in the corpus, to rise above the ocean of random combinations. In addition, the permutations would also be the basis for higher level classifications such as the INDIRECT QUESTION as PATIENT. This resulted initially in 1 120 670 feature combinations on the basis of 18411 distinct simple features, which, among others, contained 90 node-specific (morpho-syntactic), 2543 argument-specific, and 3247 extra-linguistic ones.

- (2.4) (a) `prep-and-parse-hs | prep-and-parse-sfnet-with-quotes`
- (b) `edit-fdg-context`
- (c) `merge-original-and-changes`
- (d) `post-process-hs | post-process-sfnet`
- (e) `ignore-cases`
- (f) `add-analytical-tags`
- (g) `extract-feature-combinations`
- (h) `compile-feature-table`
- (i) `select-feature-columns`
- (j) `add-feature-columns`
- (k) `find-lines-with-features-in-table`
- (l) `set-column-values-in-table`

From the general-purpose data table, only a small subset of feature columns was selected in order to keep the data manageable (e.g., all linear context data was to be excluded from this study³⁵ as well as all feature trios other than the few ones which have been used to construct selected abstract syntactic argument variables, while feature pair combinations were retained only for syntactic arguments, in addition to a minimum frequency requirement of 15 for all but the combinations of syntactic arguments and their semantic classifications) at stage (2.4i). Furthermore, the semantic classifications of non-nominal arguments (e.g., adverbs and prepositional or postpositional phrases among arguments of MANNER, QUANTITY, TIME, DURATION and FREQUENCY) were for the most part added to the table at stage (2.4j), and their context-based corrections were done at stages (2.4k-l). The resultant data table was then read as input in *R* for the following statistical analyses, which will be presented next in Section 3. A small number of variables were defined logically within *R*, namely, the general person and number features ($ANL_FIRST \leftarrow ANL_SG1 \vee ANL_PL1$; $ANL_SINGULAR \leftarrow ANL_SG1 \vee ANL_SG2 \vee ANL_SG3$, and so forth).

The final data table consisted of in all of 216 binary (logical) atomic features and 435 binary feature combinations. These broke down into 75 singular morphology-related features, 90 singular syntactic argument features (of which 22 were syntactic argument types and 68 base-form lexemes as any type of syntactic argument), 173 combinations of syntactic and semantic features, 13 combinations of syntactic and phrase-structure features, 63 combinations of syntactic argument types and base-form lexemes and 186 combinations of syntactic and morphological features, as well as 51 extralinguistic features. In addition, an ordinal index of occurrence in the research corpus, and factor (multiple-category) variables indicating the THINK lexeme, author identity and newspaper section or newsgroup as well as usage medium for each occurrence context were included in the data table. For practical purposes, the lexeme variable was supplemented with binary (logical) variables for the occurrence of each studied THINK lexeme.

³⁵ A potential avenue for further research would be to compare the analysis of simple linear context with the results of the more sophisticated syntactic argument based analysis focused on this study, that is, could one derive similar results with linear context alone, and to what extent.

3 Selection and implementation of statistical methods

3.1 Practical implementation of the employed statistical methods

For the application of statistical analysis in empirical linguistic research, Gries (2003a) has demonstrated the usefulness of a general three-tiered framework on how to proceed, consisting of 1) a univariate, 2) a bivariate, and 3) a multivariate stage,³⁶ which I will follow and to a large extent adapt but also develop further in this study (for an explicit and clear general presentation of such a framework for the purposes of exploratory data analysis, see, e.g., Hartwig and Dearing 1979: 69-79³⁷). Whereas Gries' presentation of the various applicable statistical methods is quite intertwined with his discussion of the results, I will rather first explicitly lay out and discuss the available relevant statistical methods, using example cases, and only then present the actual results in full. In accordance with Gries' example, I will begin with the analysis of all potentially interesting individual variables, that is, linguistic features, one by one, in order to identify those that are significant with respect to the linguistic phenomenon studied. In general, these variables should be anchored in earlier domain-specific research on the same subject, and have been introduced and discussed above in Section 2.2 and in depth in Appendix C.

Once this univariate analysis has identified which variables are statistically relevant and, even more importantly, which are also linguistically meaningful, I will proceed with bivariate comparisons in order to establish to what extent the individual features are associated with or dependent on each other. This may render some variables in practice redundant and, through this, it will most probably result in pruning down the number of variables for the next stage. Finally, these two stages will lead to multivariate analysis, which will further indicate the relative weights of the selected variables in relation to each other, when their joint and simultaneous influence on the studied linguistic phenomenon is taken into consideration. The crucial difference throughout between Gries' study and the one presented by me below is that whereas Gries studied a *dichotomous* alternation, my objective with the selected group of four synonymous lexemes is to extend the methodological framework to apply to the more general *polytomous* case of more than two alternatives.³⁸

In general, one should also note that the rationale for using statistical methods in this study according to a three-tiered framework is more *explorative* and *descriptive* in nature than seeking to prove pre-established hypotheses or theories (e.g., the *Processing Hypothesis* in the case of Gries 2003a). The objective here is to broaden the scope of individual contextual features and feature types which are used in the lexicographical description in general, and concerning synonymy in particular. Therefore, the confirmation of the specific results of this study will come through

³⁶ Instead of these, Gries (2003a) has used the terms *monofactorial* and *multifactorial*. As these may be confused with *Factorial Analysis (FA)* as a statistical method, I have opted for the (hopefully) less ambiguous corresponding terms *univariate* and *multivariate*, which are also the terms used by, e.g., Hartwig and Dearing (1979).

³⁷ In addition to numerical statistical analysis, Hartwig and Dearing (1979) argue forcefully for the use of visual methods in the inspection of data. In order to be able to concentrate on the former type of methods, however, I will exclude visual methods from this study.

³⁸ Alternative terms sometimes used instead of *polytomous* to refer to three or more cases are *multinomial*, *multicategory/ial* (Agresti 2002: 267) or *polychotomous* (Hosmer and Lemeshow 2000: 260).

replication, be it with other corpora (representing text types different from the ones used here), or even more preferably, with other evidence types and methods such as experimentation (cf. Arppe and Järviö 2007b), rather than from the intensive scrutiny of the significance, *Power*, *Effect Size*, or other measures in the statistical analyses. Incidentally, this attitude is aligned with recent statistical theoretical thought, represented, for instance, by Rosnow and Rosenthal (1989) and Moran (2003). Nevertheless, the range of numerical statistical analysis methods presented in *this* study will be extensive.

For undertaking all the statistical methods and the resultant analyses presented below, the public-domain *R* statistical programming environment (R Core Development Team 2007) has been used. By itself, *R* contains a vast library of already implemented, ready-to-use methods which could be applied in this study, and the number is growing fast as statisticians and researchers in other fields are contributing implementations of ever new methods and techniques.³⁹ Nevertheless, some of the methods or techniques necessary or desirable for the type of data in this study, namely, *nominal* (or *categorical*) data, and the research problem, namely, the comparison of more than two items, were not yet available in *R* when the analyses in this study were undertaken. Fortunately, the *R* environment allows the user to write functions by which he/she can implement such techniques, often building upon the functions and function libraries already existing in *R*. Therefore, some of the analyses below employ such functions written by me, which are described briefly in Appendix S. In the following presentation of the selection of methods, the function calls which provide the presented results are given at appropriate points in a distinct format, for example,

```
singular.feature.distribution(THINK.data, think.lex,  
"SX_AGE.SEM_GROUP")
```

However, understanding the full function code or the function calls requires knowledge of *R* syntax, which is beyond the scope of this dissertation; for this purpose I refer the reader to the *R* website (<http://www.r-project.org>) or textbooks introducing *R* (or its predecessors *S* and *S-PLUS*), such as Venables and Ripley (1994).

All the statistical methods use as their data a table generated from the research corpus and its subsequent analysis, using shell scripts and *R* functions as described above and in Appendix S, which is stored and publicly available in the *amph* data set at <URL: <http://www.csc.fi/english/research/software/amph>>. The main data table is designated below by the name `THINK.data`, a supplementary data table as `THINK.data.extra`, the list containing the four studied lexemes by `think.lex`, and the contextual features by labels described separately at the appropriate points.

³⁹ Especially for linguistic study, one can mention as recent contributions the `zipfR` package by Evert and Baroni (2006a, 2006b) and the `languageR` package by Baayen (2007).

3.2 Univariate methods

As Gries (2003a: 79, 107-108) points out, though linguistic phenomena are inherently influenced and determined by a multitude of variables working together at the same time, thus crying out for multivariate statistical methods, univariate analysis allows one to see in isolation the individual effects of each studied feature concerning the studied phenomenon. Such individual features are often pervasively intercorrelated so that researchers can be and have been tempted to reduce the phenomena that they study into monocausal theories, though such simple explanations are mostly inadequate (Bresnan et al. 2007; Gries 2003a: 32-36).

The singular univariate analyses below have been produced with the R function `singular.feature.distribution(data, lexemes, feature)`. For the purposes of demonstration I shall use as an example a feature which has already been studied in an earlier related study (Arppe and Järvikivi 2007b), namely, a syntactic AGENT classified as a (HUMAN) GROUP or COLLECTIVE, hereinafter denoted by the label `SX_AGE.SEM_GROUP`.⁴⁰ In order to simplify the exposition, the aggregate result of the various different univariate analyses concerning this selected feature with respect to the studied lexemes is denoted by the label `THINK.SX_AGE.SEM_GROUP`, which corresponds to the assignment of the results of a function to a variable named `THINK.SX_AGE.SEM_GROUP`, i.e.,

```
THINK.SX_AGE.SEM_GROUP <-  
singular.feature.distribution(THINK.data, think.lex,  
"SX_AGE.SEM_GROUP")
```

The starting point in univariate analysis is to compile for each studied feature a *contingency table* from the data representing the distribution of the particular feature among the studied lexemes. This can also be called a *cross-classification* or *cross-tabulation* of the studied feature and the lexemes (Agresti 2002: 36-38). In the feature-specific contingency Table 3.1 below, the frequency of the studied feature `SX_AGE.SEM_GROUP` *with* each lexeme (in the first row) is contrasted against the occurrences of each lexeme *without* the studied feature (in the second row). One should note that the features are studied here only to the extent that they occur with the selected lexemes; however often a feature may occur with other lexemes besides the selected ones, these occurrences will not be considered. This is a stance already adopted in Arppe (2002) and Arppe and Järvikivi (2007b), and it is in accordance

⁴⁰ Throughout this dissertation in terms of notation, a label in SMALL-CAPS refers to an individual linguistic feature or feature cluster, e.g., a GROUP subtype of AGENT, while a label in CAPITAL letters, e.g., `SX_AGE.SEM_GROUP`, refers to the tag by which the feature or feature cluster in question is explicitly represented in the various data tables in the `amph` data set. In the latter types of labels, the prefix `Z_XXX` refers to a (node-specific) morphological feature, `Z_ANL_XXX` to a verb-chain general morphological feature, `SX_XXX` to a syntactic argument type, `LX_XXX_YYY` to some lexeme XXX representing a part-of-speech YYY, and `SEM_XXX` to a semantic subtype of a syntactic argument or a semantic characterization of a verb chain. Furthermore, `Z_EXTRA_XXX` refers to extra-linguistic features in general, among which `Z_EXTRA_SOU_XXX` denotes one of the two main sources within the research corpus, `Z_EXTRA_DE_XXX` any of the subsections within these two sources, `Z_EXTRA_AU_XXX` the author designations in the research corpus, and `Z_EXTRA_IX_XXX` to a running identifier index assigned to each individual independent text within the two sources. Finally, `Z_PREV_XXX` refers to aspects concerning the possible repetition of the studied THINK lexemes within individual texts in the research corpus.

with *collostructional* analysis as proposed by Gries and Stefanowitsch (2004). This approach is well motivated since Gries et al. have shown that it produces results which correspond with experimental evidence, such as sentence-completion (Gries et al. 2005a) or reading times (Gries et al. 2005b), more accurately than raw counts of absolute frequencies, in which all occurrences of a feature under scrutiny are counted in.

On the basis of the raw count data derived from the corpus according to the aforementioned principle, Table 3.2 shows both feature-wise proportions (frequencies of the studied feature per lexeme relative to the overall frequency of the studied feature, in the first row) and lexeme-wise proportions (proportions of studied feature per lexeme relative to the overall frequency of lexeme, in the second row). In both Tables 3.1 and 3.2, the lexemes (i.e., columns) have been ordered according to descending feature-wise proportions (alternatively, they could be arranged in terms of descending absolute frequency per lexeme).

Table 3.1. Contingency table of the SX_AGE.SEM_GROUP feature
THINK.SX_AGE.SEM_GROUP\$ctab.ordered⁴¹

Feature/Lexeme	pohtia	harkita	miettiä	ajatella
SX_AGE.SEM_GROUP	119	64	36	37
-SX_AGE.SEM_GROUP	594	323	776	1455

Table 3.2. Feature-wise and lexeme-wise proportions of the SX_AGE.SEM_GROUP feature
THINK.SX_AGE.SEM_GROUP\$ctab.relative

Proportions/Lexeme	pohtia	harkita	miettiä	ajatella
Feature-%	46.5	25.0	14.1	14.5
Lexeme-%	16.7	16.5	4.4	2.5

Looking at the proportions in Table 3.2 we can see that almost one-half (46.5%) of the total 256 occurrences of SX_AGE.SEM_GROUP in the data are with *pohtia*; however, the proportion of SX_AGE.SEM_GROUP of all the occurrences of *pohtia* (16.7%) is practically the same as the respective proportion for *harkita* (16.5%). Furthermore, though *miettiä* and *ajatella* account for clearly lower but not negligible proportions of the overall occurrences of the SX_AGE.SEM_GROUP feature (14.1% and 14.5%, respectively), the relative proportions out of the overall occurrences of these two lexemes is substantially lower (4.4% and 2.5%, respectively). On the basis of this simple scrutiny, we could suppose that feature-wise GROUP AGENTS would appear to clearly prefer *pohtia*, but lexeme-wise both *pohtia* and *harkita* would show substantially (and equally) greater tendency for GROUP AGENTS than *miettiä* and *ajatella*.

However, we can assess and systematically construct interpretations such as these concerning the distribution represented in the contingency table with statistical means.

⁴¹ Henceforth, an identifier such as THINK.SX_AGE.SEM_GROUP\$ctab.ordered in conjunction with a Table or Figure will refer to the particular R data table, here THINK.SX_AGE.SEM_GROUP, or its subset, here ctab.ordered, which is accessed according to R syntax with the suffix \$ctab.ordered. These and other data tables are to be found in the amph data set. It is from these data tables from which the values or results represented in the Table or Figure in question have been directly derived.

As both variable types, namely, the contextual features and the lexemes, are *nominal*⁴² and *non-ordinal*⁴³ in character, the appropriate statistics concern 1) the *independence*, that is, *homogeneity*, of the distribution, and 2) the *associations* between the features and the lexemes. Both types of analysis are necessary, as they pertain to two different aspects of a relationship. For instance, a statistically significant difference in distribution might arise from the size of the sample rather than the strength of the underlying association, and likewise, very strong associations might not be supported by the significance of the distribution because the sample size may be too small.

3.2.1 Homogeneity or heterogeneity of the distribution – independence or dependence?

The first question concerns whether a studied feature is distributed evenly and uniformly among the studied lexemes or not, and what is the magnitude of the possible overall and lexeme-specific deviations from evenness. The simplest way would be to look at the absolute distribution of the feature among the studied lexemes (the first line in Table 3.1 above), in which case uniformity would entail equal absolute frequencies of the feature among the studied lexemes (i.e., with the mean frequency being naturally equal with the individual frequencies). Then, possible deviation would be evaluated as differences from the mean absolute frequency and the associated overall and individual significances. However, such one-dimensional analysis of goodness-of-fit would not really be of added informative value as it fails to take into account neither 1) the overall distributions of the studied lexemes nor 2) the distributions of other related features, for which a logically complementary distribution may hold (e.g., in the case of the SX_AGE.SEM_GROUP feature the entire set of semantic classifications for the AGENTS of the studied lexemes, or the entire set of six person-number features of FINITE verb forms).

⁴² This classification of data types into *nominal*, *ordinal*, *interval*, and *ratio* data, known as the *Stevens scale*, comes originally from Stevens (1946), who used it to prescribe what statistical methods were permissible for different types of data (Stevens 1951). As with so many things in statistics, this classification and the accompanying methodological prescriptions have been later severely criticized (see Velleman and Wilkinson 1993 for a contemporary presentation and overview of its critique), in that they are rather second-order attributes which describe how the data has been measured and to what purposes, analyses and conclusions these measures are further used, than inherent, fundamental characteristics of the data itself; however, for non-statisticians, namely linguists, the *Stevens scale* can be considered a useful guideline for selecting appropriate methods, but it is not an exception-less, absolute straightjacket.

⁴³ The lexemes could be considered to be ordered according to their overall frequency (whether according to the corpus used in this study or a general word-frequency dictionary such as FDF [Saukkonen et al. 1979]). Furthermore, the extent of their semantic fields would appear to correlate with their overall frequencies as observed in the analysis of PS (Haarala et al. 1997) above in Section 2.3.2. The most frequent of the group, *ajatella*, is described as having the broadest range of senses and usage contexts, and it is used to some extent in the descriptions of the other lexemes as a prototype, of which a semantic specification leads to the use of one or more of the three others (cf. similar to the principles of “semantic primitives” proposed by Wiezbicka 1996, Goddard 2002, and others). However, I judge these rankings as too weak and prototypicality as too abstract in order to constitute a natural, quantitative ordering which would warrant the use of ordinal methods (see Agresti 2002: 2-3; Howell 1999: 16-20; or Liebetrau 1983: 7-8, for general discussions of *nominal-ordinal* distinction), even more so as the lexemes in question were originally selected because of their roughly equal relative frequencies and the high similarity of their descriptions as well as considerable overlap of the semantic fields, in comparison to the other synonym groups scrutinized before their selection (see Section 2.1.4 above and Arppe 2002).

In an exploratory study such as this, when we want to scrutinize a large number of features, with varying degrees of logical or empirical association, I will first assess each feature on its own, in terms of its individual occurrence against its nonoccurrence among the studied lexemes, without consideration for the possible existence or frequency distribution of related complementary features, if any, among the remainders (methods for such distributions of clearly related, logically complementary features will be scrutinized and presented later in Section 3.2.3). Understood in this way, the distribution of the feature can be assessed overall with a statistical test of independence that the distribution in the contingency table deviates from the *null hypothesis* (H_0). Here, this H_0 is that the observed frequencies would equal those we could deduce and expect on the basis of the *marginal*, that is, overall feature and lexeme, frequencies (Agresti 2002: 38-39). What this null hypothesis entails in a linguistic sense is that the relative proportions of the studied feature out of the overall frequency per lexeme would be the same (even though the absolute frequencies per lexeme would vary in proportion with the overall frequencies of each lexeme), in which case the distribution would be called *homogeneous*.⁴⁴ From a linguistic viewpoint this null hypothesis represents a fully possible and conceivable state-of-affairs, rather than a *nil hypothesis* that we would *a priori* never really expect to occur at all (cf. Cohen 1994). If the null hypothesis holds, neither of the two variable types, that is, feature or lexeme, have an observable and statistically significant bearing on the other, and therefore the two variable types under scrutiny are independent of each other (in the statistical sense). In contrast, if the null hypothesis does not hold, one has reasonable grounds to assume that the *alternative hypothesis* (H_1) could be true, that is, that the two variables are dependent to some extent on each other, in which case the underlying distribution can then be considered *heterogeneous*.⁴⁵ In practice, what we evaluate is how strongly the *observed frequencies* O_{ij} represented in the contingency table deviate from the *expected frequencies* (Figure 3.1). The expected frequencies E_{ij} are calculated from the marginal row (i.e., feature) and column (i.e., lexeme) totals according to formula 3.1 below (Agresti 2002: 22, 73). The expected values for the contingency table 3.1 are shown in Table 3.3.

$$(3.1) E_{ij} = (\sum_{i=1...I} O_{ij} \cdot \sum_{j=1...J} O_{ij}) / \sum_{i=1...I} \sum_{j=1...J} O_{ij} = (R_i \cdot C_j) / N$$

where i indicates the row and j the column indexes, I indicates the number of rows and J the number of columns, R_i indicates the marginal row total of Row i and C_j the marginal column total of Column j , respectively, and N the overall total of the table.

⁴⁴ This assumption of uniformity/homogeneity is the conventional default assumption. As we will see later in Appendix K, it is possible to theoretically motivate other expectations with respect to a distribution.

⁴⁵ However, the refutation of a null hypothesis does not directly prove that the alternative hypothesis is certainly true.

Table 3.3. The expected and the marginal frequencies of the SX_AGE.SEM_GROUP feature among the studied lexemes

`chisq.test(THINK.SX_AGE.SEM_GROUP$ctab.ordered)$expected`

Feature/Lexeme	pohtia	harkita	miettä	ajatella	$\sum_{\text{row}} = R_i$
SX_AGE.SEM_GROUP	53.6	29.1	61.1	112.2	256
-SX_AGE.SEM_GROUP	659.4	357.9	750.9	1379.8	3148
$\sum_{\text{column}} = C_j$	713	387	812	1492	3404

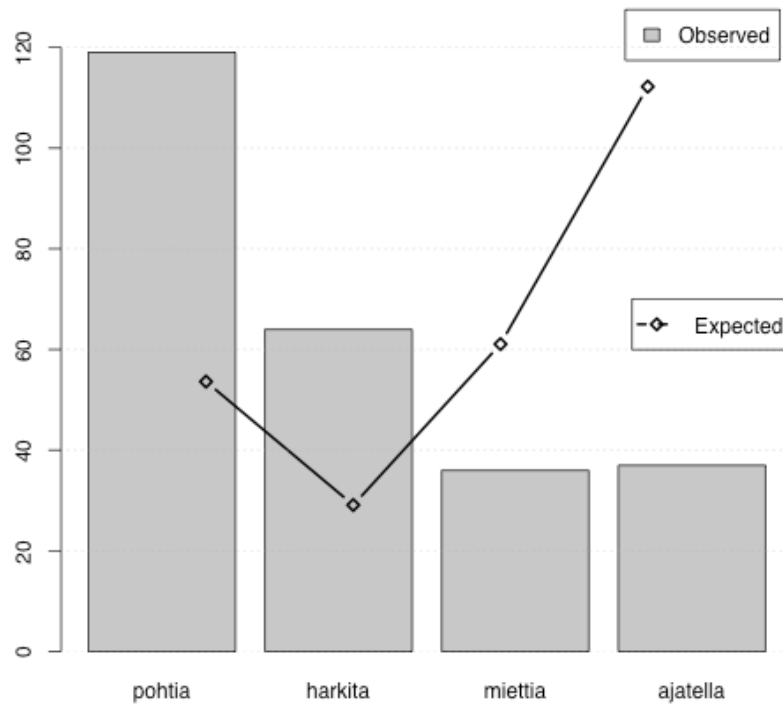


Figure 3.1. The observed and the expected frequencies for the SX_AGE.SEM_GROUP feature among the studied lexemes (in accordance with the test of independence).

The overall homogeneity, or the lack of it, that is, heterogeneity, in the contingency table can be assessed with two commonly used *approximate*⁴⁶ measures, namely, the *Pearson chi-squared* (X^2) or the *Likelihood-ratio chi-squared* (G^2) statistics (e.g., Agresti 2002: 78-80). Whereas both the X^2 and the G^2 statistics asymptotically converge with each other as well as the theoretical chi-squared (χ^2) distribution they approximate as the overall frequency increases, with smaller overall frequencies their behavior differs to some extent. However, for very small overall frequencies both methods are unreliable in the judgement of independence. In the last-mentioned case, small-sample methods such as *Fisher's exact test* can be used instead (Agresti 2002: 91-94; Pedersen 1996), but their use has been historically limited due to their extremely high computational cost, and often the theoretical considerations of the

⁴⁶ *Approximate* here means that we use a formula which is an approximation of what (statistical) theory would dictate, but which is relatively simple in mathematical terms and computationally inexpensive. N.B. Often, we only have approximate formulas at our disposal, as in many cases the underlying theoretical equations do not have exact numerical solutions.

scientific field in question render questionable the study of low-frequency phenomena. For both X^2 and G^2 , simulation and other studies have provided indication of minimum requirements in order to retain sufficient reliability (e.g., the so-called *Cochran* [1952, 1954] conditions; see also Agresti [2002: 8, 395-396]; or the *minimum average expected frequency* by Roscoe and Byars [1971]), which take into account the size of the contingency table and the expected values in the individual cells of the table. For a 2x4 table as is the case in the univariate studies here, and taking into consideration the overall frequencies of the studied lexemes, the minimum overall frequency for individual features is 24 occurrences.⁴⁷ Since my focus is, in the spirit of Sinclair et al. (1991, 2001: ix-x), on the more common contextual associations of the studied lexemes, those which are generally established in the linguistic community rather than the less frequent ones, be they exceptions, idiolectal preferences, or mere random linguistic variation, I content myself with studying features equaling and exceeding this minimum value, as this cut-off point will also substantially prune the overall number of features. Nevertheless, with an overall frequency of 256 occurrences the SX_AGE.SEM_GROUP feature certainly clearly exceeds both minimum frequencies. Of the two approximate measures, the Pearson statistic is somewhat simpler to calculate and the interim components of its calculation can directly be used in its follow-up analysis, so I will use it in the subsequent analysis (as did Gries 2003a: 80-83). The formula for calculating the X^2 is given in 3.2 (Agresti 2002: 78-79, formula 3.10):

$$(3.2) X^2 = \sum_{i=1}^I \sum_{j=1}^J [(O_{ij} - E_{ij})^2 / E_{ij}]$$

Where i and j are the row and column indices, I and J the number of rows and columns, respectively, and N the overall total.

This overall X^2 value together with the appropriate *degrees of freedom* is then used to yield an estimate of the *level of significance* according to the chi-squared (χ^2) distribution.⁴⁸ In general, the number of degrees of freedom for the X^2 (as well as G^2)

⁴⁷ Cochran's assessment (presented in its first form in 1952: 334, and with more specifications in 1954: 420) is that a minimum expected value ≈ 1 in some cells is acceptable as long as at least 80% of the other expected values are > 5 . For the 2x4 contingency tables (adding up to 8 cells altogether) used in the univariate analyses here – assuming that the overall frequencies of the individual features are substantially less than the overall frequencies of the studied lexemes, so that we can focus on the expected values in the feature frequency row – this entails that at most one (or with a stretch two, as 80% of 8 = 1.6 cells) of the expected values for features per lexeme can be around 1 while the other three have to be ≥ 5 . With the overall lexeme frequencies fixed, the minimum overall feature frequency which satisfies these conditions is 24 ($\sum(\text{FEATURE})=24 \Rightarrow E(\text{ajatella}|\text{FEATURE})=24 \cdot 1492/3404=10.52$; $E(\text{miettiä}|\text{FEATURE})=5.73$; $E(\text{pohtia}|\text{FEATURE})=5.03$; $E(\text{harkita}|\text{FEATURE})=2.73$), or 21 in the less conservative two-cell interpretation ($\sum=21 \Rightarrow E(\text{ajatella}|\text{FEATURE})=21 \cdot 1492/3404=9.2$; $E(\text{miettiä}|\text{FEATURE})=5.01$; $E(\text{pohtia}|\text{FEATURE})=4.4$; $E(\text{harkita}|\text{FEATURE})=2.39$). Roscoe and Byars (1971) argue that we should rather assess the *average expected frequency* than the individual minimum expected frequencies; for this they propose values ranging from 2 to 10. In any case, adhering to Cochran's conditions always entails conformance with Roscoe and Byars's slightly more lenient minimum requirements.

⁴⁸ Modern statistical programs such as *R* yield directly an exact P-value (e.g., `chisq.test(THINK.SX_AGE.SEM_GROUP$ctab.ordered)$p.value` or `THINK.SX_AGE.SEM_GROUP$omnibus.p`). Without such applications, the X^2 value calculated from the contingency table can be compared manually in a table with a pre-calculated value corresponding to the χ^2 with the appropriate degrees of freedom and the predetermined critical P-value (α), often denoted as $\chi^2(\alpha, df)$ (which in this case, with the critical significance level $P < .05$ and $df=3$,

statistic is $df=(I-1)\cdot(J-1)$, where I is the number of rows and J the number of columns of the contingency table, being in this case $df=(2-1)\cdot(4-1)=1\cdot3=3$. The significance level, often also known as the *P-value* or *alpha* (α), indicates the probability that the observed values in the contingency table could have been sampled by chance from the assumed underlying population. In the human sciences, to which linguistics certainly belongs, along with psychology, sociology, and other disciplines, the *critical P-value* or *critical alpha*, that is required for an observation to be considered statistically significant, is conventionally (and N.B. quite arbitrarily) set at $P<0.05$, and this P-value will also be used throughout in this study (e.g., Howell 1999: 128-129).⁴⁹ This particular critical P-value entails that there is a 5% risk (or chance, whichever way one sees it) that the observations in question could have been sampled from the population by chance. In other words, that 1 in 20 sampled observations with this particular P-level are more likely to be the results of random sampling variation, than representative of a real lack of independence in the assumed underlying population of which the observations are a sample. One should make note that this type of significance testing assesses the probability of how likely it is that we would observe our data, given the null hypothesis of independence, i.e., $P(DATA|H_0)$, and *not* vice versa, i.e., $P(H_0|DATA)$. Furthermore, rejecting the null hypothesis on the account of a significant P-value does not in a complementary sense amount to a direct confirmation of the alternative hypothesis, which is typically the actually sought-after conclusion (Cohen 1994).

Looking at formula 3.2, we can see that it consists of cell-by-cell calculations of the squared deviations of the observed values from the expected values, normalized by the expected values. These cell-by-cell calculations are known as X^2 *contributions* and their square roots as *Pearson residuals* (Agresti 2002: 81, formula 3.12). In order to calculate the overall X^2 value for the contingency table, we can first calculate the X^2 contributions, shown in Table 3.4. The sum of the X^2 contributions, and thus the overall $X^2=197.07$ ⁵⁰, is substantially more than the critical value, i.e., $\chi^2(\alpha=0.05, df=3)=7.81$, and the corresponding exact $P(X^2=197.07, df=3)=1.81e^{-42}$ is clearly below the critical value. This indicates that we can reject the null hypothesis of independence between the feature and the lexemes, and assume (though not definitely conclude) instead that there is a strong association between the type of lexeme and the particular feature.

would be achieved by `qchisq(alpha=.05,df=3,lower.tail=FALSE)` or `THINK.SX_AGE.SEM_GROUP$omnibus.min`). In such a case, if $X^2 > \chi^2(\alpha, df)$, the observed distribution in the contingency table is considered statistically significant.

⁴⁹ In exploratory analysis as is the case in this study, Cohen (1992: 156) in fact makes the suggestion that a critical P-value as high as $\alpha<.10$ could be acceptable, but I will nevertheless (mostly) adhere to the established convention.

⁵⁰ The more precise X^2 value provided by *R* is actually =197.0691; however, for reasons of readability I will present rounded values throughout this dissertation, though the underlying actual calculations are naturally carried out with unrounded values.

Table 3.4. X^2 contributions of the SX_AGE.SEM_GROUP feature among the studied lexemes, with a sign supplemented to signify whether the observed value exceeded (+) or subceeded (-) the expected value.

`chisq.test(THINK.SX_AGE.SEM_GROUP$ctab.ordered)$residuals^2, or THINK.SX_AGE.SEM_GROUP$cellwise["cell.stat"]` for only the feature-specific values we are actually interested in.

Feature/Lexeme	pohtia	harkita	miettiä	ajatella
SX_AGE.SEM_GROUP	+79.71	+41.84	-10.29	-50.41
(-SX_AGE.SEM_GROUP)	(-6.48)	(-3.40)	(-0.84)	(-4.10)

Assessing Power and Effect Size

This now standard practice in statistical analysis of focusing on testing the significance of a null hypothesis, and the associated focus/fixation on a dichotomous decision to reject, or not to reject, the null hypothesis on the basis of some particular pre-selected P-value has been criticized by, for example, Rosnow and Rosenthal (1989) and Cohen (1990, 1992, 1994). In their view, with apparent justification, this practice has led to the widespread neglect of the three other relevant statistical variables, namely, *Power* ($1-\beta$), minimum sample size (N), and *Effect Size* (ES , denoted in the case chi-square test of independence as w). Instead, they rather recommend a combined consideration of these variables together with the significance level (α), preferably in order to establish the minimum sample size necessary on the basis of the three other criteria. Alternatively, they recommend the assessment of the *Effect Size* and/or *Power* afterwards in addition to reporting the significance level. Specifically, highly significant P-values should not be interpreted as automatically reflecting large effects (Rosenthal and Rosnow 1989: 1279).

For contingency tables of the size studied here, with $df=3$, fixing $\alpha=0.05$ and *Power* ($1-\beta$) at 0.20 (i.e., $\beta=0.80$), Cohen (1988: 258, Table 7.4.6, or 1992, Tables 1 and 2) has calculated as the minimum sample sizes 1090 for a *small effect* ($w=0.10$), 121 for a *medium effect* ($w=0.30$), and 44 for a *large effect* ($w=0.50$). These aforementioned three designations of *Effect Sizes* are generic conventions proposed by Cohen, which can be used in the case that they cannot be estimated from prior research or otherwise. As the overall sample in this study far exceeds Cohen's highest minimum sample size for detecting small effects ($3404 > 1090$), we can assume that the amount of data is sufficient for discovering even quite small effects.⁵¹

Nevertheless, heeding this critique and advice, I will calculate *post hoc* the *Effect Sizes* as well as the *Power* of the individual univariate analyses. The formula for the *Effect Size* w (Cohen 1988: 216-221, formula 7.2.1, or Cohen 1992: 157) for a chi-squared test derived from a contingency table and the associated formula for *Power* ($1-\beta$) (following Agresti 2002: 243-244, formulas 6.8 and 6.9), together with interim calculations of the *noncentrality* parameter λ and the probability of *Type II errors* (β) are given in 3.3-3.6 below. As one can see, formula 3.3 structurally resembles X^2 statistic, with relative proportions (=probabilities) instead of absolute observed and expected frequencies; furthermore, the degrees of freedom are the same for all the

⁵¹ See Cohen (1998: 253-267, Tables 7.4.1-7.4.115) for minimum sample sizes N with a range of other values of α , *Power*, and *Effect Size* w than the ones presented here.

formulas, that is, for the contingency tables of the size studied here $df=(4-1)\cdot(2-1)=3$. In fact, we will note in Section 3.2.2 below that Effect Size w is closely related to measures of association based on the chi-squared statistic, and that w can be calculated from those measures. In particular, for the 2x4 tables scrutinized here – or generally speaking any table with either two rows or two columns, as $\min(2,J)=\min(I,2)=2$ – Effect Size is equal to Cramér's V .

$$(3.3) w = \left\{ \sum_{i=1}^I \sum_{j=1}^J [P(O_{ij}) - P(E_{ij})]^2 / P(E_{ij}) \right\}^{1/2}$$

$$= \left\{ \sum_{i=1}^I \sum_{j=1}^J [O_{ij}/N - E_{ij}/N]^2 / (E_{ij}/N) \right\}^{1/2}$$

so that $\sum_{i=1}^I \sum_{j=1}^J [P(O_{ij})] = 1$ and $\sum_{i=1}^I \sum_{j=1}^J [P(E_{ij})] = 1$

$$(3.4) \lambda = N \cdot w^2$$

$$(3.5) \beta = P[X^2_{df,\lambda} > \chi^2(df, \alpha)], \quad df = (I-1) \cdot (J-1)$$

$$(3.6) \text{Power} = 1 - \beta$$

Where i and j and the row and column indices, I and J and the number of rows and columns, respectively, and N the overall total.

For the purpose of transparency, the probabilities of the observed and the expected frequencies, as well as the cell-wise contributions to the Effect Size statistic w_{ij} , with respect to the SX_AGE.SEM_GROUP feature among the studied lexemes are presented in Tables 3.5-3.7 below. For instance, for the co-occurrence of the SX_AGE.SEM_GROUP feature with the lexeme *pohtia*, the probability of observed occurrence, corresponding to the alternative hypothesis H_1 , is $P(O_{ij}) = O_{ij}/N = 119/3404 = 0.0350$,⁵² and the probability of expected occurrence, corresponding to the null hypothesis H_0 , is $P(E_{ij}) = E_{ij}/N = 53.62/3404 = 0.0158$. Thus, the cell-wise contribution $w_{SX_AGE,SEM_GROUP,pohtia} = (0.0350 - 0.0158)^2 / 0.0158 = 0.0234$. Consequently, the Effect Size is the square root of the sum of the individual cell values w_{ij} , which is $w = (0.0579)^{1/2} = 0.241$. Moving further, the noncentrality parameter $\lambda = 3404 \cdot 0.241^2 = 197.07$, and $\beta = P[X^2_{df=3, \lambda=197.0691} > \chi^2(df=3, \alpha=0.05)] = P[\chi^2(df=3, \alpha=0.05), \lambda=197.07, df=3] = P[7.815, \lambda=197.07, df=3] = 2.43e^{-30} \approx 0$, finally yielding Power as $1 - 0 = 1.0$.

Table 3.5. Probabilities of the observed frequencies of the SX_AGE.SEM_GROUP feature among the studied lexemes.

Probabilities of Observed frequencies/Lexeme	pohtia	harkita	miettiä	ajatella
SX_AGE.SEM_GROUP	0.0350	0.0188	0.0106	0.0109
-SX_AGE.SEM_GROUP	0.175	0.0950	0.228	0.427

⁵² Here it is my understanding that the alternative hypothesis H_1 is fixed to equal the state of affairs represented by the observed distribution, though such an H_1 is, of course, only one of the many possible distributions which would deviate significantly from the homogeneous distribution corresponding to the null hypothesis.

Table 3.6. Probabilities of the expected frequencies of the SX_AGE.SEM_GROUP feature among the studied lexemes.

Probabilities of Observed frequencies/Lexeme	pohtia	harkita	miettiä	ajatella
SX_AGE.SEM_GROUP	0.0158	0.00855	0.0179	0.0330
-SX_AGE.SEM_GROUP	0.194	0.105	0.221	0.405

Table 3.7. Cell-wise contributions to the *Effect Size* statistic *w*.

Cell-wise contributions	pohtia	harkita	miettiä	ajatella
SX_AGE.SEM_GROUP	0.0234	0.0123	0.00302	0.0148
-SX_AGE.SEM_GROUP	0.00190	0.0001000	0.000246	0.00120

Cell-wise follow-up scrutiny – identifying where and how strong the deviations are

Notwithstanding the above critique, the X^2 test (or any other test of significance) by itself tells us whether there is something very significant overall somewhere in the relationship between the studied feature and lexemes, as is certainly the case for Table 3.1. However, the X^2 test says very little about the exact locus or the direction of this association. Statisticians have long urged researchers to supplement tests of significance with studies concerning the nature of the association (Agresti 2002: 80). Probably the simplest method is to study cell-by-cell the X^2 contributions, shown above in Table 3.4. Firstly, we can scrutinize to what extent individual cells account for the overall deviation from the expected values. A conservative procedure is to assess whether some individual cells by themselves exceed the minimum value required by the critical P-value (α) with the same degrees of freedom as the entire table, or is the overall X^2 value actually the sum of smaller deviations. A less conservative procedure is to regard each of the individual cells as their own tables, having thus $df=1$ and consequently a lower minimum critical X^2 statistic value. Secondly, for those cell-wise contributions that we do deem significant, we can look in which direction, either above or below, the observed values lie in relation to the expected values.

According to the conservative procedure, compared against the minimum X^2 value for the entire table $\chi^2(df=3, \alpha=0.05)=7.81$, we can see in Table 3.4 above that the X^2 contributions of all the feature-specific cells clearly exceed this value. When we then take into consideration the direction of the observed values in comparison to the expected values, we can conclude that both *pohtia* and *harkita* have been observed significantly more with the SX_AGE.SEM_GROUP feature than would be expected if this feature occurred evenly, whereas the case is the contrary for both *miettiä* and *ajatella*. The same results naturally hold when the X^2 contributions are compared to the minimum single-cell X^2 value $\chi^2(df=1, \alpha=0.05)=3.84$ in accordance with the less conservative procedure. A further step would be to use the exact P-values associated with the cell-wise X^2 contributions to quantify the significance of the deviations, as these can easily be calculated with the help of most modern statistical programs (in our case with the function call `THINK.SX_AGE.SEM_GROUP$cellwise["cell.p"]` in R). For the conservative procedure, the P-values are $3.54e^{-17}$ for *pohtia*, $4.34e^{-09}$ for *harkita*, $1.63e^{-02}$ for

miettiinä, and $6.54e^{-11}$ for *ajatella*. As can be seen, the significance of the deviation for *miettiinä* is considerably less than for the three other verbs.

Another closely related method which conveniently combines the assessment of the significance and direction of the cell-by-cell contributions is to calculate the *standardized Pearson residuals*, for which the formula is 3.7 (Agresti 2002: 81, formula 3.13). In the relatively small contingency table that we are now studying, a standardized Pearson residual which exceeds at least 2 in absolute value indicates a significant deviation in the cell in question. For larger tables the minimum absolute value should be 3 or even more, but no exact values have been provided in the literature. So, for the cell with the count for the co-occurrence of the SX_AGE.SEM_GROUP feature with the lexeme *pohtia*, the standardized Pearson residual is $(119-53.62)/[53.62 \cdot (1-256/3404) \cdot (1-713/3404)]^{1/2} = +10.44$. For the rest of the studied lexemes, the respective values are +7.14 for *harkita*, -3.82 for *miettiinä*, and -9.85 for *ajatella* (which we can obtain with the function call `THINK.SX_AGE.SEM_GROUP$cellwise["residual.pearson.std"]`). All of these values clearly exceed 2 in absolute terms (or 3, for that matter), so all the cell-wise deviations can be considered significant. From the signs of the respective values we can make the conclusions that the SX_AGE.SEM_GROUP feature occurs in conjunction with both *pohtia* and *harkita* significantly more than expected, and with both *miettiinä* and *ajatella* significantly less than expected. These are exactly the same results that we obtained by studying the X^2 contributions.

$$(3.7) e_{ij/\text{standardized Pearson residual}} = (O_{ij} - E_{ij}) / [E_{ij} \cdot (1 - R_i/N) \cdot (1 - C_j/N)]^{1/2}$$

Where i and j and the row and column indices, I and J and the number of rows and columns, R_i and C_j are the row and column marginal totals, respectively, and N the overall total.

A third way to assess the components of the distribution of the SX_AGE.SEM_GROUP feature among the studied lexemes is to conduct pairwise comparisons, selecting iteratively the appropriate lexeme columns for the calculation of simple 2x2 Pearson chi-squared tests. This is similar to the study of *contrasts* in the *Analysis of Variance*, applicable for interval data. As can be seen from the results shown below in Table 3.8, there are significant differences between the paired comparisons of all the verbs except *pohtia* and *harkita*. This could be linguistically interpreted as stratifying the studied lexemes into three groups, with *pohtia* and *harkita* forming a single group, and *miettiinä* and *ajatella* each forming a group of their own.

Table 3.8. Pairwise comparisons of the SX_AGE.SEM_GROUP feature among the studied lexemes: P-values of pairwise X^2 tests, with significant tests marked with (*)
`THINK.SX_AGE.SEM_GROUP$pairwise["pair.p"]`

Lexeme/Lexeme	pohtia	harkita	miettiinä	ajatella
pohtia	-	$9.84e^{-01}$	$*5.36e^{-15}$	$*1.28e^{-33}$
harkita	$9.84e^{-01}$	-	$*3.05e^{-12}$	$*3.44e^{-27}$
miettiinä	$*5.36e^{-15}$	$*3.05e^{-12}$	-	$*1.50e^{-02}$
ajatella	$*1.28e^{-33}$	$*3.44e^{-27}$	$*1.50e^{-02}$	-

The problem with such pairwise comparisons is that in the case of a relatively small group of items (say, less than 5 as is the case here) they can, in principle, stratify too much or too little. This may be the case if the comparisons of each immediately

adjacent, frequency-wise descending pairing are statistically significant, or if these adjacent pairings are nonsignificant, even when the overall distribution and some longer-distance pairing(s) may be significant. In terms of interpretation, the pairwise comparisons can only establish a gradient of greater to lesser association of the individual lexemes with respect to the studied feature, as the overall benchmark (in the form of the expected values) derivable from the entire distribution is explicitly not used. Therefore, at least in the case of relatively small group of semantically closely related lexemes such as we have here, the follow-up measures concerning the cell-wise contributions or their variants are more attractive, and simpler, too, and I will subsequently focus on them in the presentation of the results. The overall behavior of the three different methods presented above in the assessment of cell-wise contributions with respect to the entire range of studied features, namely, 1) the comparison of cell-wise contributions against the minimum X^2 value with the same df as the entire table, 2) the comparison of cell-wise contributions against the minimum X^2 with $df=1$, and 3) the standardized Pearson residuals, will be presented later in Section 4.1.1 covering the results. In order to ease the analysis *en masse* of a large number of singular features as is the case in this study, the results of these various cell-wise strategies can be simplified according to whether individual cells do, or do not, indicate a significant deviation from the expected distribution, and in which direction the deviation lies in relation to the expected distribution (see Table 3.9).

Table 3.9. Simplified representation of the various methods of assessing cell-wise contributions for the distribution of the SX_AGE.SEM_GROUP feature among the studied lexemes, with (+) denoting a significant observed deviation above the expected value, (–) a significant observed deviation below the expected value, and (0) a nonsignificant observed deviation.

Assessment strategy	Minimum significant value	pohtia	harkita	miettiinä	ajatella
Table minimum	$\chi^2(df=3, \alpha=0.05) > 7.81$	+	+	–	–
Cell-wise minimum	$\chi^2(df=1, \alpha=0.05) > 3.84$	+	+	–	–
Standardized Pearson residual	$ e_{ij}/\text{standardized Pearson residual} > 2$	+	+	–	–

Adjusting the critical P-levels in follow-up analyses

For such follow-up, that is, *post hoc*, analyses, it has been traditional in many scientific fields, though not in all fields and not consistently, to require adjusted lower critical P-values for such tests to be considered significant (for a relevant example in linguistics, see Gries 2003a: 81-82). The oldest and simplest such procedure is known as the *Bonferroni correction*, which has been followed by many modifications and alternatives. The rationale behind these adjustment procedures lies in the risk/chance of encountering a randomly significant distribution that the critical P-value (α) represents. Once we have established (for a contingency table with more than 2 rows and 2 columns) that the entire distribution is statistically significant with a given pre-selected critical P-level, if we then, after the fact, decide (or even if we already planned this beforehand) to continue with a large number of pairwise or other follow-up analyses of the individual contributions concerning the same contingency table, we run the risk, in principle, of encountering just such a false significance by chance. In order to retain the so-called *family-wise error rate*, which is the aggregate probability that *at least one* in the family/set of follow-up tests is nonsignificant, as equal to the

overall critical P-level, there exists an overabundance of different procedures to choose from, each which calculate a set of adjusted critical P-values, comparing them with the actual P-values obtained with the follow-up tests, often by considering the distribution of the follow-up P-values as a whole. Each of these procedures emphasizes a different aspect, controlling either *Type I errors* (α), that is, the probability of mistakenly classifying nonsignificant effects as significant (represented by the traditional simple Bonferroni procedure of dividing the critical P-level by the number of *post hoc* tests, $\alpha_{Bonferroni} = \alpha_{familywise}/n$, or less conservative sequential methods such as Hochberg 1988; Holland and Copenhaver 1988; Hommel 1988; Rom 1990), or the *false discovery rate*, that is, the probability of nonsignificant effects being correctly rejected as nonsignificant, leading to a better *Power* ($1-\beta$) (Benjamini and Yekutieli 2001). In a comparison of a range of α -controlling procedures, Olejnik et al. (1997) judged the Rom (1990) procedure to have the greatest *Power*.

Lately, this practice has been severely criticized by scholars from a variety of fields (Perneger 1998; Moran 2003; O'Keefe 2003; Nakagawa 2004), mainly because it can drastically reduce the Power of experiments to reveal interesting effects. In addition, there exists no formal consensus, or, in the opinion of some, there is a fundamental inconsistency, regarding the situations when it should be applied, that is, what exactly constitutes a family of tests for which the family-wise error rate should be controlled. Perhaps the most convincing argument is that as more and more research is conducted, spurious results are inevitable and thus in effect ultimately uncontrollable; however, such results will be falsified, that is, they will *not be reproduced* by later research. To the contrary, it is extremely improbable that *all* results will be spurious, even when *some* certainly will turn out to be so (Moran 2003: 405). For example, the fact that all the X^2 contributions in addition to the overall X^2 statistic concerning Table 3.1 above are highly significant strongly supports for the conclusion that the observations concerning SX_AGE.SEM_GROUP are (for the most part) truly significant. The consequent alternative approach, then, is to report and judge the P-values as they are, and this is the practice chosen in this study, too. As the number of variables exceeding the minimum frequency threshold is quite high at 477, this will naturally entail that some of the judgements of significance may probably be incorrect, in the order of 20-30 ($\approx 477/20$).

A Zipfian alternative perspective

As an alternative complement to the above analysis, we could also assess univariate feature distributions in terms of how they relate to *Zipf's law*. However, the set of four lexemes that I am scrutinizing in this dissertation is yet a very small number of elements for the application of Zipfian principles, which typically concern the entire lexeme inventory of a corpus, or at the very least all the members of a given synonym set. Nevertheless, I will in Appendix K explore ways of scrutinizing distributions of features, for even such a small set as the selected THINK lexemes, from a Zipfian perspective.

3.2.2 Measures of association

The statistical analysis presented thus far has first and foremost concerned whether the observed distribution incorporates an effect or a relationship of some type that can be considered statistically significant, the exact nature of which has been left unspecified. Consequently, the second question we can pose regarding an observed distribution focuses on the characteristics, direction, and strength of the relationship between the studied variables. In the case of nominal variables, such a relationship is in general referred to by the term *association*, instead of *correlation* which is reserved for association between interval (and often also rank-ordered) variables. For the measurement of association, there is available a wide selection of different methods, of which some, typically the older “traditional” ones, are based on the chi-squared test and in effect attempt to normalize its overall result with respect to the size of data in various ways. Other methods evaluate the extent to which knowing one (independent) variable would allow one to predict or determine the (dependent) other variable according to differing premises, understood generally in terms of *Proportionate Reduction of Error*, or alternatively, *Proportion of Explained Variation* (often referred to with the acronym PRE). As all of these methods attempt to summarize the relationship between the two variables over and above all the individual cell-wise comparisons, they are called *summary measures*. Since many of the nominal methods are applicable only to dichotomous variables with 2x2 tables, rather than polytomous variables as is the case in this study (with always more than two lexemes, sometimes more than two contextual features), the number of relevant methods presented below is conveniently pruned.

Cramér’s V

The association measures based on the chi-squared statistic X^2 , which are applicable for polytomous (nominal) variables, are 1) Pearson’s *Contingency Coefficient* (or *Coefficient of mean square contingency*) C (Goodman and Kruskal 1954: 739, formula 5; Liebetrau 1983: 13, formula 3.1; Garson 2007),⁵³ 2) Tschuprow’s *Contingency Coefficient* T (Goodman and Kruskal 1954: 739-740, formula 6; Liebetrau 1983: 14, formula 3.3; Garson 2007), and 3) Cramér’s V (Cramér 1946: 282-283, 443-444; see also Goodman and Kruskal 1954: 740, formula 7; Liebetrau 1983: 14-15, formula 3.4; Garson 2007). Of these three methods, Cramér’s V is considered the best measure of association because of its norming properties, in that it ranges between 0–1 and can in practice always attain either end-point values regardless of the dimensions of the table (Goodman and Kruskal 1954: 740; see also Buchanan 1974: 643). Therefore, it is one of the measures selected and used in this study. The formula for Cramér’s V is given below in formula 3.8; its significance level is equal to that of the underlying X^2 statistic. For instance, for all the 2x4 singular feature tables in this study, $q=\min(4,2)=2$, and $N=3404$, so as the X^2 statistic for the SX_AGE.SEM_GROUP feature is 197.07, Cramér’s V is consequently $V=\{197.07/[\overline{3404 \cdot (2-1)}]\}^{1/2}=0.241$ and the associated P-value is $P(X^2=197.07, df=3)=1.81e^{-42}$. As was noted above, Cramér’s V is closely linked to the estimation of Effect Size w and the associated *Power* for chi-squared tests, with the relationship presented in formula 3.9. This gives rise to the notion that measures of

⁵³ This measure has alternatively been referred to as ϕ by Liebetrau (1983).

association in general could be considered to indirectly estimate *Effect Size*. Like all chi-squared based measures, Cramér's V is symmetrical, so it provides us with a single and simply computable value by which we can rank the studied individual features in terms of their associations. Such symmetric statistics are often considered the nominal equivalents of the well-known Pearson's correlation coefficient r for interval data.

(3.8) $V = \{X^2/[N \cdot (q-1)]\}^{1/2}$, where $q = \min(I, J)$, i.e., the lesser of the table dimensions I and J , N the overall total, and X^2 calculated according to formula 3.2 above.⁵⁴

(3.9) $w = \{\sum_{i=1}^I \sum_{j=1}^J [O_{ij}/N - E_{ij}/N]^2 / (E_{ij}/N)\}^{1/2} = (X^2/N)^{1/2}$
 So, $w/(d-1)^{1/2} = \{X^2/[N \cdot (d-1)]\}^{1/2} = V$; and therefore $w = V \cdot (d-1)^{1/2}$

Measures based on Proportionate Reduction of Error (PRE)

The disadvantage of Cramér's V , together with all the other chi-squared based measures, is that they are connected with the underlying distribution and dimensions of the contingency table determined by the number of classes in the polytomous variables. Therefore, the values of these measures are meaningfully comparable only when the overall table frequencies and dimensions are the same (Goodman and Kruskal 1954: 740). Consequently, we can with justification compare the values of Cramér's V for the 2x4 singular feature contingency tables scrutinized in this study, but it would not be meaningful to compare these values with the respective ones of another study with, say, five lexemes instead, with some other overall lexeme frequencies. Due to this inherent lack of universal comparability, the Proportionate Reduction of Error (PRE) measures are an attractive alternative and supplement to the chi-squared based tests. What PRE measures in principle evaluate is how much the proportion of classification errors can be reduced (e.g., Costner 1965), or alternatively how much more of the variation of the dependent variable can be explained and accounted for (e.g., Kviz 1981), when knowing some aspect of the distribution of the dependent variable conditional on the independent variable, in comparison to some baseline knowledge. The latter is typically understood as knowing a given aspect of the overall distribution of the dependent variable (see general formula 3.10 below, which applies for all PRE methods, Reynolds 1977: 32-34; also Agresti 2002: 56-57). Probably the most commonly known and widely used asymmetric PRE methods applicable for polytomous (nominal) data are 1) the Goodman-Kruskal *lambda* ($\lambda_{A|B}$ and $\lambda_{B|A}$), 2) the Goodman-Kruskal *tau* ($\tau_{A|B}$ and $\tau_{B|A}$), and 3) Theil's *Uncertainty Coefficient* (UC or Theil's $U_{A|B}$ and $U_{B|A}$). Of these measures, the Goodman-Kruskal λ has been used earlier in a similar linguistic study (Gries 2003a: 126).

(3.10) Proportionate Reduction of Error (PRE) = $[P_{\text{Error/baseline}} - P_{\text{Error/measure}}] / P_{\text{Error/baseline}}$

⁵⁴ Interestingly, Cramér himself does not appear to give the symbolic designation for this statistic attributed to himself, but rather presents a way of norming Pearson's coefficient (referred by him also as its square ϕ^2) so that the values will always fall between [0,1]; neither does he explicitly suggest presenting a square root of this normed measure. Where the latter convention originates from is unclear to me, as for instance Goodman and Kruskal (1954: 740) present the measure in a nameless form.

Goodman-Kruskal λ

The asymmetric Goodman-Kruskal $\lambda_{A|B}$ was originally conceptually proposed by Guttman in 1941, but named and promoted by Goodman and Kruskal (1954: 740-747, formulas 9-10; see also Liebetrau 1983: 17-24, formulas 3.12, 3.13, 3.15 and 3.16; Agresti 2002: 69; Garson 2007). This statistic $\lambda_{A|B}$ can be interpreted as how much knowing both the independent variable B and the maximum of the corresponding dependent variable A conditional on B [i.e. $\max(A|B)$] increases our chances of correctly predicting A , compared to a baseline of knowing only the overall distribution and the maximum of the dependent variable A [i.e. $\max(A)$]. The opposite case of $\lambda_{B|A}$ is the same except that variables are interchanged so that the independent variable is A and the dependent variable is B . The formula for both versions of the asymmetric Goodman-Kruskal $\lambda_{A|B}$, with A denoting the Column variable and B the Row variable, are given in 3.11 and 3.12. The $\lambda_{A|B}$ statistic is well-defined, provided that not all (non-zero) occurrences are crammed into one row, or into one column in the case of the $\lambda_{B|A}$ statistic. Alternatively put, this requirement means that at least two rows, or at the least two columns, respectively, must each have at least one non-zero cell (Liebetrau 1983: 19).

$$(3.11) \lambda_{\text{Row|Column}} = [\sum_{k=1 \dots I} \max(O_{i=k,j}) - \max(R_i)] / [N - \max(R_i)]$$

$$(3.12) \lambda_{\text{Column|Row}} = [\sum_{k=1 \dots J} \max(O_{i,j=k}) - \max(C_j)] / [N - \max(C_j)]$$

where i and j are the row and column indices, I and J the number of rows and columns, R_i and C_j the row and column marginal totals, respectively, and N the overall total.

Thus, in the case of the `SX_AGE.SEM_GROUP` feature and the studied lexemes, with *Feature* as the Row variable and *Lexeme* as the Column variable, prior knowledge of the lexeme and each lexeme's individual distribution with respect to the occurrence and nonoccurrence of the `SX_AGE.SEM_GROUP` feature increases our understanding beyond the baseline knowledge of the overall distribution of `SX_AGE.SEM_GROUP` feature with

$$\lambda_{\text{Feature|Lexeme}} = \{[\max(O_{\text{Feature},\text{pohtia}}) + \max(O_{\text{Feature},\text{harkita}}) + \max(O_{\text{Feature},\text{mieltiä}}) + \max(O_{\text{Feature},\text{ajatella}})] - \max(R_{\text{SX_AGE.SEM_GROUP}}, R_{\text{-SX_AGE.SEM_GROUP}})\} / [N - \max(R_{\text{SX_AGE.SEM_GROUP}}, R_{\text{-SX_AGE.SEM_GROUP}})]$$

$$= [(594+324+776+1455) - 3148] / (3404 - 1492) = 0.00052.$$

Likewise, prior knowledge of the feature's occurrence (or its nonoccurrence) and the corresponding lexeme distributions, compared to the baseline of knowing only the overall distributions of the studied lexemes, yields

$$\lambda_{\text{Lexeme|Feature}} = \{[\max(O_{\text{SX_AGE.SEM_GROUP},\text{Lexeme}}) + \max(O_{\text{-SX_AGE.SEM_GROUP},\text{Lexeme}})] - \max(C_{\text{pohtia}}, C_{\text{harkita}}, C_{\text{mieltiä}}, C_{\text{ajatella}})\} / [N - \max(C_{\text{pohtia}}, C_{\text{harkita}}, C_{\text{mieltiä}}, C_{\text{ajatella}})]$$

$$= [(119+1455) - 1492] / (3404 - 1492) = 0.0429.$$

The relevant cell-values in the calculations of $\lambda_{\text{Feature|Lexeme}}$ and $\lambda_{\text{Lexeme|Feature}}$ have been highlighted below in Tables 3.10 and 3.11, respectively.

Table 3.10. Relevant cell values for the calculation of $\lambda_{Feature|Lexeme}$, with the selected maxima $\max(O_{Feature,pohtia}, \dots)$ and $\max(R_{SX_AGE.SEM_GROUP}, R_{-SX_AGE.SEM_GROUP})$ in boldface.

Feature/Lexeme	pohtia	harkita	miettä	ajatella	$\sum_{row}=R_i$
SX_AGE.SEM_GROUP	119	64	36	37	256
-SX_AGE.SEM_GROUP	594	323	776	1455	3148
$\sum_{column}=C_j$	713	387	812	1492	3404

Table 3.11. Relevant cell values for the calculation of $\lambda_{Lexeme|Feature}$, with the selected maxima $\max(C_{pohtia}, C_{harkita}, C_{miettä}, C_{ajatella})$ and $\max(O_{SX_AGE.SEM_GROUP, Lexeme})$ in boldface.

Feature/Lexeme	pohtia	harkita	miettä	ajatella	$\sum_{row}=R_i$
SX_AGE.SEM_GROUP	119	64	36	37	256
-SX_AGE.SEM_GROUP	594	323	776	1455	3148
$\sum_{column}=C_j$	713	387	812	1492	3404

Goodman-Kruskal τ

The asymmetric Goodman-Kruskal $\tau_{A|B}$ ⁵⁵ was originally suggested by W. Allen Wallis, but was formulated explicitly by Goodman and Kruskal (1954: 745-747, formula 17; see also Liebetrau 1983: 24-31, formulas 3.24, 3.25, 3.27 and 3.28). This statistic $\tau_{A|B}$ is analogous to $\lambda_{A|B}$, but the focus is on the prediction of expected probabilities of all the classes of the dependent variable rather than the discrete choices of only one of its classes at a time. Therefore, $\tau_{A|B}$ can be interpreted as how much knowing both the independent variable B and the overall distribution of the dependent variable A conditional on B (i.e., $A|B$) increases our accuracy in predicting the *probabilities* of (all) the various classes of A [i.e. $P(A|B)$], compared to a baseline of knowing only the overall probabilities of the classes of the dependent variable A [i.e. $P(A)$]. In a gambling analogy, the baseline for the Goodman-Kruskal $\lambda_{A|B}$ corresponds to the chance of success when betting always only on the most frequent dependent item B for each independent A , without any consideration for the outcome history. In contrast, the baseline for the Goodman-Kruskal $\tau_{A|B}$ reflects the chance of success in betting in the long run, while knowing the entire expected distribution of the dependent B for each independent A , and keeping track of accumulating outcomes. Here, too, the calculation of $\tau_{B|A}$ is the same except that variables are interchanged so that the independent variable is A and the dependent variable is B . The formulas for both versions of the asymmetric Goodman-Kruskal tau $\tau_{A|B}$, with A denoting the Column variable and B the Row variable, are given in 3.13 and 3.14. The $\tau_{A|B}$ statistic is well-defined if at least two cells are non-zero and these cells are in separate rows; in the case of the $\tau_{B|A}$ statistic these two non-zero cells have to be in separate columns (Liebetrau 1983: 26).

$$(3.13) \tau_{Column|Row} = [N \cdot \sum_{i=1 \dots I} \sum_{j=1 \dots J} (O_{ij}^2 / R_i) - \sum_{j=1 \dots J} (C_j^2)] / [N^2 - \sum_{j=1 \dots J} (C_j^2)]$$

$$(3.14) \tau_{Row|Column} = [N \cdot \sum_{i=1 \dots I} \sum_{j=1 \dots J} (O_{ij}^2 / C_j) - \sum_{i=1 \dots I} (R_i^2)] / [N^2 - \sum_{i=1 \dots I} (R_i^2)]$$

where i and j are the row and column indices, I and J the number of rows and columns, R_i and C_j the row and column marginal totals, respectively, and N the overall total.

⁵⁵ Goodman and Kruskal (1954) actually refer to this statistic as the λ_* , some others as the *lambda-max* λ_{max} .

Thus, in the case of the SX_AGE.SEM_GROUP feature and the studied lexemes, with *Feature* as the Row variable and *Lexeme* as the Column variable, prior knowledge of the lexemes and their individual distributions with respect to the occurrence and nonoccurrence of the SX_AGE.SEM_GROUP feature increases our understanding beyond the baseline knowledge of the overall distribution of SX_AGE.SEM_GROUP feature with

$$\begin{aligned} \tau_{Feature|Lexeme} &= \{N \cdot [(O_{SX_AGE.SEM_GROUP,pohtia}^2/C_{pohtia}) + (O_{-SX_AGE.SEM_GROUP,pohtia}^2/C_{pohtia}) + \dots \\ &+ (O_{SX_AGE.SEM_GROUP,ajatella}^2/C_{ajatella}) + (O_{-SX_AGE.SEM_GROUP,ajatella}^2/C_{ajatella})] - \\ &(R_{SX_AGE.SEM_GROUP}^2 + R_{-SX_AGE.SEM_GROUP}^2)\} / [N^2 \cdot (R_{SX_AGE.SEM_GROUP}^2 + R_{-SX_AGE.SEM_GROUP}^2)] \\ &= \{3404 \cdot [(119^2 + 594^2)/713 + (64^2 + 323^2)/387 + (36^2 + 776^2)/812 + (37^2 + 1455^2)/1492] - \\ &(256^2 + 3148^2)\} / [3404^2 - (256^2 + 3148^2)] = 0.0579. \end{aligned}$$

Likewise, prior knowledge of the feature's occurrence (or its nonoccurrence) and the corresponding lexeme distributions, compared to the baseline of knowing only the overall distributions of the studied lexemes, yields

$$\begin{aligned} \tau_{Lexeme|Feature} &= \{N \cdot [(O_{SX_AGE.SEM_GROUP,pohtia}^2/R_{SX_AGE.SEM_GROUP}) + \dots \\ &+ (O_{SX_AGE.SEM_GROUP,ajatella}^2/R_{SX_AGE.SEM_GROUP}) + (O_{-SX_AGE.SEM_GROUP,pohtia}^2/R_{-SX_AGE.SEM_GROUP}) \\ &+ \dots + (O_{-SX_AGE.SEM_GROUP,ajatella}^2/R_{-SX_AGE.SEM_GROUP})] - (C_{pohtia}^2 + C_{harkita}^2 + C_{mieltia}^2 + C_{ajatella}^2)\} / \\ &[N^2 \cdot (C_{pohtia}^2 + C_{harkita}^2 + C_{mieltia}^2 + C_{ajatella}^2)] \\ &= \{3404 \cdot [(119^2 + 64^2 + 36^2 + 37^2)/256 + (594^2 + 323^2 + 776^2 + 1455^2)/3148] - \\ &(713^2 + 387^2 + 812^2 + 1492^2)\} / [3404^2 - (713^2 + 387^2 + 812^2 + 1492^2)] = 0.0211. \end{aligned}$$

As is evident, in comparison to λ all cell-values are equally relevant in the calculations of $\tau_{Feature|Lexeme}$ and $\tau_{Lexeme|Feature}$.

Theil's uncertainty coefficient

Theil's uncertainty coefficient U (Theil 1970: 125-129, formula 13.6; see also Agresti 2002: 57, formula 2.13; Garson 2007) is similar to the Goodman-Kruskal τ in that it also takes into consideration the entire expected distribution of the dependent variable for each independent variable. The difference is that U is based on the concept of *entropy* from information theory rather than the estimated probability of occurrence, and the statistic calculates the reduction of entropy rather than that of prediction error. Here, entropy is understood to represent the average *uncertainty* concerning the value of the dependent variable, when knowing the determining independent variable. However, the two approaches are interconnected in that entropy is defined as (minus) the expected value of the logarithm of the probability (Theil 1970: 127). The formula for the Uncertainty Coefficient is given in 3.15 and 3.16 below. The uncertainty coefficient U is well-defined even in the case that some cells have zero occurrences, increasing thus its attractiveness (remembering that $\lim_{x \rightarrow 0} [x \cdot \log(x)] = 0$, see Theil 1970: 127).

$$(3.15) U_{Row|Column} = [H(X) + H(Y) - H(XY)] / H[X]$$

$$(3.16) U_{Column|Row} = [H(Y) + H(X) - H(XY)] / H[Y]$$

Where

$$H(X) = -\sum_{i=1...I} [(R_i/N) \cdot \log_e(R_i/N)];$$

$$H(Y) = -\sum_{j=1...J} [(C_j/N) \cdot \log_e(C_j/N)]; \text{ and}$$

$$H(XY) = -\sum_{i=1...I} \sum_{j=1...J} [(O_{ij}/N) \cdot \log_e(O_{ij}/N)],$$

and i and j are the row and column indices, I and J and the number of rows and columns, R_i and C_j the row and column marginal totals, respectively, and N the overall total.

Thus, in the case of the `SX_AGE.SEM_GROUP` feature and the studied lexemes, with *Feature* as the Row variable and *Lexeme* as the Column variable, the row-specific (horizontal) frequency-wise entropy is

$$H(X) = -[(R_{SX_AGE.SEM_GROUP}/N) \cdot \log_e(R_{SX_AGE.SEM_GROUP}/N) + (R_{-SX_AGE.SEM_GROUP}/N) \cdot \log_e(R_{-SX_AGE.SEM_GROUP}/N)] = -[(256/3404) \cdot \log_e(256/3404) + (3148/3404) \cdot \log_e(3148/3404)] = 0.266;$$

the column-specific (vertical) feature-wise entropy is

$$H(Y) = -[(C_{pohtia}/N) \cdot \log_e(C_{pohtia}/N) + \dots + (C_{ajatella}/N) \cdot \log_e(C_{ajatella}/N)] = -[(713/3404) \cdot \log_e(713/3404) + (387/3404) \cdot \log_e(387/3404) + (812/3404) \cdot \log_e(812/3404) + (1492/3404) \cdot \log_e(1492/3404)] = 1.278;$$

and the joint entropy is

$$H(XY) = -[(O_{SX_AGE.SEM_GROUP,pohtia}/N) \cdot \log_e(O_{SX_AGE.SEM_GROUP,pohtia}/N) + (O_{-SX_AGE.SEM_GROUP,pohtia}/N) \cdot \log_e(O_{-SX_AGE.SEM_GROUP,pohtia}/N) + \dots] = -[(119/3404) \cdot \log_e(119/3404) + (594/3404) \cdot \log_e(594/3404) + \dots] = 1.518.$$

Therefore, in terms of entropy, prior knowledge of the lexeme and each lexeme's individual distribution with respect to the occurrence and nonoccurrence of the `SX_AGE.SEM_GROUP` feature increases our understanding beyond the baseline knowledge of the overall distribution of `SX_AGE.SEM_GROUP` feature with $U_{Lexeme|Feature} = (0.267 + 1.278 - 1.518) / 1.278 = 0.0213$. Likewise, prior knowledge of the feature's occurrence (or its nonoccurrence) and the corresponding lexeme distributions, compared to the baseline of knowing only the overall distributions of the studied lexemes, yields $U_{Lexeme|Feature} = (0.267 + 1.278 - 1.518) / 0.267 = 0.102$.

Significance of association measure values

In general, measures of association can be calculated for data of any size, as these statistics do not make any assumptions concerning some hypothesized underlying population, but rather try to interpret and represent the data as it is. However, significance values can in principle be estimated also for the PRE measures presented here, provided that the sample size is sufficiently large. One could well wonder what meaning these P-values would have on top of the chi-squared based evaluation of whether a frequency distribution incorporates a statistically significant relationship. Basically, the significance values for PRE measures of association provide an estimate of how probable the observed, calculated value of the measure in question is in comparison to a hypothetical (zero) value, given the marginal values of the particular sampled distribution. The formulas for the variances of the various

measures are even more convoluted than the calculations for the measures themselves, and are therefore not presented in this dissertation, though they have been implemented by me in the R functions described briefly in Appendix S. The P-values for the PRE measures calculated above concerning the distribution of the SX_AGE.SEM_GROUP feature among the studied lexemes, both lexeme-wise and feature-wise, are presented in Table 3.12. As can be seen, all measures of association presented above are clearly significant with respect to the relation between the studied feature and lexemes. Further on, I will provide significance values for the association measures only occasionally.

Table 3.12. The statistics (\hat{E}) and significance values (P-values) of the selected symmetric and asymmetric nominal measures of association, calculated both lexeme-wise, i.e., $\hat{E}(Feature|Lexeme)$, and feature-wise, i.e., $\hat{E}(Lexeme|Feature)$, for the relationship between the SX_AGE.SEM_GROUP feature and the studied lexemes.

THINK.SX_AGE.SEM_GROUP\$associations

Measure	$\hat{E}(Feature Lexeme)$	P(\hat{E})	$\hat{E}(Lexeme Feature)$	P(\hat{E})
V	0.241	1.81e ⁻⁴²	=	=
$\lambda_{A B}$	0	NA	0.0429	1.93e ⁻¹¹
$\tau_{A B}$	0.0579	1.87e ⁻⁴²	0.0211	2.53e ⁻⁴⁶
$U_{A B}$	0.102	5.66e ⁻⁴⁰	0.0213	5.66e ⁻⁴⁰

Comparison of the characteristics of the present association measures

The various summary measures of association differ along several parameters according to which they can be classified (Weisberg 1974; see also Buchanan 1974; Garson 1975: 200-202; Liebetrau 1983: 85-88; and Garson 2007). Therefore, it is recommended that association measures be selected according to the fit of these parameters with the general characteristics of the studied phenomenon and the focus of the particular research question. However, the comparison of the different measures also indicates that no single method is perfect, as they fundamentally differ in the theoretical concepts on which they are based, and thus they ultimately measure different things. Consequently, it is recommended that researchers use more than one method and capitalize on the differences that they possibly bring out, preferably selecting methods which vary with respect to their underlying assumptions. This is motivated by Weisberg's observation that methods based on the same conceptual model correlate to a great degree (Weisberg 1974: 1639, 1647-1648, 1652; see also Reynolds 1977: 50). Additional, practical factors are the ease of computation, or lately, whether a particular method has been implemented in the available statistical software or not. Furthermore, a pragmatic factor to consider is whether to include methods that have been used earlier in similar studies in the scientific field in question, in order to achieve some level of comparability and continuity with earlier research. However, prior use is no automatic guarantee of appropriateness⁵⁶; (cf. Buchanan 1974: 625-626, Weisberg 1974: 1646).

⁵⁶ In addition to the Goodman-Kruskal λ used by Gries (2003a: 126, Note 5 to Chapter 6) to study a similar linguistic research question, he also applied the Somer's d (2003: 82) and the r^2 (a variant expression of Pearson's r) measures (Gries 2003a: 126, Note 8 to Chapter 6). Of these, Somer's d (see, e.g., Liebetrau 1983: 77-82, formulas 5.52a and 5.52b; Agresti 2002: 68; Garson 2007) requires ordinal data and Pearson's correlation coefficient r (see, e.g., Liebetrau 1983: 45-49, formula 4.9) interval data, to which types none of the variables scrutinized in this study belong.

In any case, a researcher should at the least be aware of the conceptual basis (and subsequent implications regarding their interpretation) of the measures he/she has selected and used. These are presented and discussed at length in Appendix L. A concise summary of the theoretical properties of all selected association measures presented above are provided below in Table 3.13. Summary comparisons of their values and their correlations and covariation for the range of features scrutinized in this study will be covered later in Section 4.1.1, the presentation of the general results.

Table 3.13. Theoretical properties of selected association measures applicable for polytomous nominal data (adapted from Weinberg 1974, Charts 6–7, with supplements from Garson 1975, 2007, and my own observations presented here)

Measure	Perfect relationship	Null relationship	Causal directionality	Sensitivity to marginals	Inter-mediate values
V (Cramér)	Moderate	Independence	Symmetric	Sensitive	Linear, non-smooth
λ (Goodman-Kruskal Lambda)	Moderate (predictive)	Accord	Asymmetric	Sensitive	Linear, non-smooth
τ (Goodman-Kruskal Tau)	Moderate (predictive)	Independence	Asymmetric	Sensitive	Curvilinear
U (Theil's Uncertainty Coefficient)	Moderate (predictive)	independence	asymmetric	Sensitive	Curvilinear

Verbal characterization of association measure values

Finally, various verbal characterizations based on differing threshold values have been suggested for interpreting nominal association measures, specifically those belonging to the PRE type, and the underlying relationship. For instance, Corbett and Le Roy (2002: 189) suggest designating relationships with PRE association values in the range 0.0–0.10 as *very weak*, 0.10–0.19 as *weak*, 0.20–0.29 as *moderate*, and 0.30–1.0 as *strong*. Towards the more rigorous end, Kviz (1981: 419) tentatively suggests the value of $\hat{E}=0.5$ as a cutting for PRE measures point, with higher values consequently representing a *strong* relationship and lower values a *weak* relationship, on the grounds that this would conceptually anchor the interpretation in terms of whether a *majority* of the variance in the relationship is explained or not.

On the other hand, Howell (1999: 186) notes that even relatively low association values may in practice represent noteworthy relationships, especially when a large number of factors is involved in the studied phenomenon. This is typically the case in the complexities of real human behavior, or when even a relatively small improvement resulting from a better understanding of the phenomenon is valuable (e.g., reduction of mortality from accident, injury or disease in human societies). As a matter of fact, the univariate results will yield just such seemingly low association values.

Furthermore, although the aforementioned threshold figures may appear relatively insubstantial, we must remember that in the case of PRE measures they indicate *added* explanatory power, over and above some default levels based on the frequency

of the most common outcome or the overall frequencies of all of the possible outcomes. Nevertheless, we should bear in mind that the aforementioned verbal interpretations and the associated threshold values are all more or less arbitrary, though they do provide generally applied reference points handy in pruning and reducing data and variables when no clear, natural divisions are evident.

3.2.3 Grouped univariate analysis for a set of related contextual features

The statistical methods presented hitherto have focused on the distribution of a single contextual feature among the studied lexemes. This has been expressed in terms of the occurrence or nonoccurrence of the feature in question, where the nonoccurrences can include some other, logically related and possibly complementary features. In fact, it is possible to scrutinize at the same time, using the same methods presented above, groups of such related features, interpreting these as different categories or classes of the same variable. For instance, human GROUPS and COLLECTIVES are not the only semantic type of AGENT that the studied lexemes can have. Quite obviously, human INDIVIDUALS (denoted henceforth by the label SX_AGE.SEM_INDIVIDUAL) are another and even more frequent type ($n=2251$) of AGENT; in fact, the corpus analysis demonstrates that there are in all 9 different possible semantic classification of the AGENTS for the studied lexemes. However, 7 of these 9 observed semantic types have very low relative frequencies, namely, abstract NOTIONS ($n=7$), EVENTS involving people ($n=5$), physical ARTIFACTS ($n=4$), FAUNA ($n=2$), ACTIVITY ($n=2$), manifestations of COMMUNICATION ($n=2$), and LOCATIONS ($n=1$), and thus fall below any threshold for meaningful statistical analysis; in addition, there are seven unclassified miscellaneous instances of AGENTS.⁵⁷ So, instead of contrasting the observed occurrences of the SX_AGE.SEM_GROUP feature against its nonoccurrences among the studied lexemes, we can study its distribution against the other related and frequent semantic classification SX_AGE.SEM_INDIVIDUAL. The corresponding contingency table containing the observed frequencies of the two studied features among the four lexemes is presented in Table 3.14, and the corresponding feature-wise and lexeme-wise relative proportions in Tables 3.15 and 3.16, respectively. These have been calculated using the R function `multiple.feature.distribution`, i.e.,

```
THINK.SX_AGE.SEM_INDIVIDUAL_GROUP
<- multiple.feature.distribution(THINK.data, think.lex,
c("SX_AGE.SEM_INDIVIDUAL", "SX_AGE.SEM_GROUP"))
```

Table 3.14. Contingency table representing the observed frequencies of the related two features of AGENT subtypes among the studied lexemes.

```
THINK.SX_AGE.SEM_INDIVIDUAL_GROUP$ctab.ordered
```

Feature/Lexeme	ajatella	miettä	pohtia	harkita	$\Sigma(\text{Feature})$
SX_AGE.SEM_INDIVIDUAL	1047	632	374	198	2251
SX_AGE.SEM_GROUP	37	36	119	64	256
$\Sigma(\text{Lexeme})$	1084	668	493	262	2507

⁵⁷ Many of these other semantic classifications can, in fact, be understood as manifestations of human groups, i.e., LOCATIONS used to refer to a group of people living or working there and ACTIVITIES and EVENTS used to refer to recurrent or one-time congregations of groups of people for some particular purpose. The remaining classifications, i.e., abstract NOTIONS, elements of COMMUNICATION and ARTIFACTS, refer to anthropomorphic uses.

Table 3.15. Lexeme-wise proportions of the related two features relative to the overall frequencies of the studied lexemes.

THINK.SX AGE.SEM INDIVIDUAL GROUP\$ctab.relative.lexeme

Feature/Lexeme (%)	ajatella	miittiä	pohtia	harkita	Σ (Feature)
SX_AGE.SEM_INDIVIDUAL	70.2	77.8	52.5	51.2	62.9
SX_AGE.SEM_GROUP	2.5	4.4	16.7	16.5	10.0
Σ(Lexeme)	72.7	82.2	69.2	67.7	72.9

Table 3.16. Proportions of the studied lexemes relative to the overall frequencies of each of the two related features.

THINK.SX AGE.SEM INDIVIDUAL GROUP\$ctab.relative.feature

Feature/Lexeme (%)	ajatella	miittiä	pohtia	harkita	Σ (Feature)
SX_AGE.SEM_INDIVIDUAL	46.5	28.1	16.6	8.8	100.0
SX_AGE.SEM_GROUP	14.5	14.1	46.5	25.0	100.0

By simply looking at the raw counts and corresponding proportions presented in the Tables 3.14-3.16 above, we can observe clear tendencies. In general, either one of the two main semantic classes of AGENTS occurs with a large majority of each studied lexeme (67.7–82.2%). However, observed as proportions of overall frequencies of the lexemes, both *ajatella* and *miittiä* have clearly larger proportions of INDIVIDUAL agents, 72.7% and 82.2%, respectively, than the other two lexemes. In contrast, both *pohtia* and *harkita* have clearly larger proportions of GROUP AGENTS, 16.7% and 16.5%, respectively, than the two other lexemes. In terms of proportions out of the overall feature frequencies, *ajatella* and *miittiä* account for the majority of occurrences of INDIVIDUAL AGENTS, whereas *pohtia* and *harkita* do the same for GROUP AGENTS. On the basis of these figures, I can propose the hypothesis that while *ajatella* and *miittiä* are associated with INDIVIDUAL AGENTS, *pohtia* and *harkita* are associated with GROUP AGENTS.

As was demonstrated above, such observations can be systematically evaluated and confirmed with the application of statistical methods. Firstly, we can test the homogeneity of this observed distribution with the chi-squared (X^2) test of independence between the two variables, that is, the studied four lexemes and the two related features. The cell-wise contributions to the X^2 statistic are presented in Table 3.17, summing up to 233.62, which with $df=(4-1)\cdot(3-1)=3$ clearly exceeds the critical minimum of $\chi^2(\alpha=0.05, df=3)=7.815$ and is significant with $P(X^2=233.62, df=3)=2.28e^{-50}$. Therefore, I can conclude that the two variables, comprised of the features on the one hand and the lexemes on the other hand, are not at all independent of each other. This is supported in that the *Effect Size* for the observed distribution is $w=0.305$, amounting to a *medium* effect according to Cohen's proposed benchmarks, and the associated *Power* is very strong with $(\beta-1)=1.0$. Furthermore, it is worth noting that this observed effect for the combination of the two semantic types of AGENT is somewhat higher than that which was observed earlier for the SX_AGE.SEM_GROUP feature alone (where $w=0.241$).

Table 3.17. Chi-squared (X^2) contributions for the related two features among the studied lexemes; all cells are statistically significant cells (with $df=3$).

THINK.SX AGE.SEM INDIVIDUAL GROUP\$cell.stat

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX_AGE.SEM_INDIVIDUAL	5.58	1.73	10.65	5.90
SX_AGE.SEM_GROUP	49.06	15.21	93.64	51.85

We can then look for the foci of the divergences from the independent, homogeneous distribution with cell-wise analysis of X^2 contributions. Calculated conservatively against the overall $df=3$, some, but not all cells exceed the minimum value, yielding corresponding P-values presented in Table 3.18. On the basis of this analysis, *pohtia* can be judged to be negatively associated with INDIVIDUAL AGENTS, whereas the other three lexemes would appear neutral with respect to this semantic AGENT type. The contrast is clearer with GROUP AGENTS, where *ajatella* and *miettiä* are significantly negatively associated and *pohtia* and *harkita* are significantly positively associated with this AGENT type. When we compare these results against the standardized Pearson residuals presented in Table 3.19, we can see that this latter method is again less conservative, since all cells clearly exceed the minimum threshold values (being either $e_{ij}>2$ or $e_{ij}<-2$). Indeed, if instead of the conservative threshold with $df=3$ we compare the cell-wise X^2 contributions with the most lenient critical value $\chi^2(\alpha=.05, df=1)=3.841$, all cells except one exceed this value, the sole exception being *miettiä* in conjunction with an INDIVIDUAL AGENT. These results concur with the pairwise comparison of *miettiä* and *pohtia* by Arppe and Järvikivi (2007b) with respect to the GROUP AGENTS. For INDIVIDUAL AGENTS, however, this four-lexeme comparison distinguishes *pohtia* from the rest as dispreferring INDIVIDUAL AGENTS, which was not found in the earlier study. This difference may result from the inclusion of all person/number features as INDIVIDUAL agents here, in comparison to the scrutiny of only the FIRST and THIRD PERSON SINGULAR in the earlier study.

Table 3.18. Significance values of the chi-squared (X^2) contributions for the related two features of AGENT subtypes among the studied lexemes; statistically significant cells (with $df=3$) in boldface.

THINK.SX AGE.SEM INDIVIDUAL GROUP\$cell.p

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX_AGE.SEM_INDIVIDUAL	1.34e ⁻⁰¹	0.630	- 1.38e⁻⁰²	1.17e ⁻⁰¹
SX_AGE.SEM_GROUP	- 1.27e⁻¹⁰	- 0.00164	+ 3.62e⁻²⁰	+ 3.22e⁻¹¹

Table 3.19. Standardized Pearson residuals for the related two features of AGENT subtypes among the studied lexemes; all cells are significant cells, i.e., $|e_{ij}|>2$.

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX_AGE.SEM_INDIVIDUAL	+9.81	+4.81	-11.39	-8.03
SX_AGE.SEM_GROUP	-9.81	-4.81	+11.39	+8.03

In addition to the assessment of the homogeneity of a distribution, we can also calculate the various summary measures of associations between the two features and the studied lexemes, which are presented in Table 3.20. We can see that the symmetric Cramér's V is equal to the Effect Size w (as we are dealing with a $2 \times N$ table, where $q=\min[2, N]=2$). Interpreting Cramér's V in terms of the explained variance of the studied lexemes, the observed association of the studied lexemes and their major types of AGENTS is not insignificant, as should be expected. Furthermore, the lexeme-wise asymmetric association measures treating the features as predictable

dependents are all higher than the opposite-direction associations; in addition, all of these measures are significant. Accordingly, in terms of interpretation, knowing the lexeme alone can be understood to account for approximately one-tenth of the behavior of the studied lexemes with respect to their occurrence with the two semantic types of AGENTS.

Table 3.20. The statistics (\hat{E}) and significance values (P-values) of the selected nominal measures of association, calculated both lexeme-wise, i.e., $\hat{E}(Feature|Lexeme)$, and feature-wise, i.e., $\hat{E}(Lexeme|Feature)$, for relationship between the related SX_AGE.SEM_INDIVIDUAL and SX_AGE.SEM_GROUP features and the four studied lexemes.

THINK.SX_AGE.SEM_INDIVIDUAL_GROUP\$associations

Measure	$\hat{E}(Feature Lexeme)$	$P(\hat{E})$	$\hat{E}(Lexeme Feature)$	$P(\hat{E})$
V	0.305	$2.28e^{-50}$	=	=
$\lambda_{A B}$	0	NA	0.0576	$1.35e^{-11}$
$\tau_{A B}$	0.0932	$2.39e^{-50}$	0.0315	$4.65e^{-51}$
$U_{A B}$	0.129	$4.17e^{-46}$	0.0336	$4.176e^{-46}$

More features and larger tables

The number of related singular features to be scrutinized at the same time need not be restricted to only two alternatives as was the case above. For instance, we can extend our earlier study (Arppe and Järvikivi 2007b, see also Appendix K) of the FIRST PERSON SINGULAR feature (Z_SG1) to cover all person/number morphological features observable in the data, that is, the SECOND and THIRD PERSONS SINGULAR and the FIRST, SECOND, and THIRD PERSONS PLURAL (denoted by the corresponding labels Z_SG2, Z_SG3, Z_PL1, Z_PL2, and Z_PL3). The contingency table presenting the observed occurrences of all these person/number features with respect to the studied lexemes is presented in Table 3.21. The relative divisions of these features among the lexemes as well as the relative proportions of these features of the overall frequencies of the studied lexemes are presented in Tables 3.22 and 3.23, respectively. At first glance, we can see that some features are considerably rarer in the observed corpus than others, namely, all the FIRST and SECOND PERSON PLURAL features. Furthermore, certain features seem to account for larger proportions of some lexemes than is the case for others, for instance, the FIRST PERSON SINGULAR with *ajatella* and *miettiinä*, the SECOND PERSON SINGULAR with *miettiinä* (and *ajatella*), the THIRD PERSON SINGULAR with *pohitia*, and the THIRD PERSON PLURAL with *ajatella*. Overall for such a large table, assessing the raw count data and considering all their comparisons is more difficult than the cases presented earlier.

Table 3.21. Contingency table presenting the frequencies of the occurrences of the six related person/number features among the studied lexemes.

THINK.Z PERSON NUMBER\$ctab.ordered

Feature/Lexeme	ajatella	miettiä	pohtia	harkita	Σ(Feature)
Z_SG1	170	57	9	12	248
Z_SG2	93	73	3	2	171
Z_SG3	163	126	177	43	509
Z_PL1	14	4	0	3	21
Z_PL2	17	17	15	2	51
Z_PL3	91	21	37	15	164
Σ(Lexeme)	548	298	241	77	1164

Table 3.22. Lexeme-wise proportions of the six related person/number features relative to the overall frequencies the studied lexemes.

THINK.Z PERSON NUMBER\$ctab.relative.lexeme

Feature/Lexeme (%)	ajatella	miettiä	pohtia	harkita	×(Feature)
Z_SG1	11.4	7.0	1.3	3.1	5.7
Z_SG2	6.2	9.0	0.4	0.5	4.0
Z_SG3	10.9	15.5	24.8	11.1	15.6
Z_PL1	0.9	0.5	0.0	0.8	0.6
Z_PL2	1.1	2.1	2.1	0.5	1.5
Z_PL3	6.1	2.6	5.2	3.9	4.5
Σ(Lexeme)	36.6	36.7	33.8	19.9	31.8

Table 3.23. Feature-wise proportions of the studied lexemes relative to the overall frequencies of each of the six related person/number features.

THINK.Z PERSON NUMBER\$ctab.relative.feature

Feature/Lexeme (%)	ajatella	miettiä	pohtia	harkita	Σ(Feature)
Z_SG1	68.5	23.0	3.6	4.8	100.0
Z_SG2	54.4	42.7	1.8	1.2	100.0
Z_SG3	32.0	24.8	34.8	8.4	100.0
Z_PL1	66.7	19.0	0.0	14.3	100.0
Z_PL2	33.3	33.3	29.4	3.9	100.0
Z_PL3	55.5	12.8	22.6	9.1	100.0

Again, we can test the overall relationship between the four lexemes, on the one hand, and the six related features, on the other, with the test of the homogeneity of the distribution in the entire table. The cell-wise contributions to the chi-squared (X^2) statistic are given in Table 3.24, and sum up to $X^2=224.48$, which for this 6x4 table with a subsequent $df=(6-1) \cdot (4-1)=15$ also clearly exceeds the minimum value of $\chi^2(\alpha=0.05, df=15)=25.00$ and is highly significant with $P(224.48, df=115)=2.16e^{-39}$. Furthermore, the Effect Size is $w=0.439$ with a maximum corresponding $Power=1.0$. On the basis of all these figures we can conclude that overall the studied lexemes and six person-number features are interrelated.

Next, we want to know where in particular the foci of the detected overall divergence are located. Compared against the conservative minimum statistic value 25.00, with $df=15$, only three cells by themselves exceed this value. All of these are with *pohtia*, which occurs significantly less with the FIRST and SECOND PERSON SINGULAR features, but significantly more with THIRD PERSON SINGULAR feature. These are reflected quite naturally also in the cell-wise P-values in Table 3.25, Firstly, this indicates that the overall divergence arises from many relatively smaller deviations, but secondly also

that with a larger table and consequently higher degrees of freedom, as is the case here, the conservative cell-wise assessment may become too stringent. Indeed, when we look at the standardized Pearson residuals in Table 3.26, we can see that a larger proportion of individual cells exceed the critical minimum value (either $e_{ij} < -2$ or $e_{ij} > +2$). Now, *ajatella* is positively associated with both FIRST and SECOND PERSON SINGULAR and THIRD PERSON PLURAL features, but negatively associated with THIRD PERSON SINGULAR and SECOND PERSON PLURAL features. Furthermore, *miettiinä* is positively associated with the SECOND PERSON SINGULAR and negatively with the THIRD PERSON PLURAL features, whereas *pohtia* is positively associated with the THIRD PERSON SINGULAR and negatively with the FIRST and SECOND PERSON SINGULAR and FIRST PERSON PLURAL features, while *harkita* is positively associated with THIRD PERSON SINGULAR feature and negatively with the FIRST PERSON PLURAL feature. Looking at the associations from the feature-wise perspective, the SECOND PERSON SINGULAR seems the most discriminatory, with *ajatella* and *miettiinä* associated positively and *pohtia* and *harkita* negatively with it; similar but less sweeping deviations can be noted for all the other features, too.

However, a sizable proportion of cells remain below the critical level when studied as standardized Pearson residuals, thus retaining this less conservative strategy as a discriminatory tool. Therefore, in light of the overall cell-wise assessment results presented earlier for singular feature analysis and here for grouped-feature cell-wise analysis, I find the use of the standardized Pearson residuals as the most attractive strategy (see Table 3.27 for a comparison of the results in simplified form according to the notation presented earlier in conjunction with the singular feature analysis). Finally, in comparison to the earlier pairwise comparison of *miettiinä* and *pohtia* by Arppe and Järvikivi (2007b), the results for this four-lexeme scrutiny are quite similar with respect to the FIRST PERSON SINGULAR, with *pohtia* being negatively associated and *miettiinä* neutral with this feature. However, we must remember that these corpus-based results were shown in this earlier study not to represent the whole truth concerning the semantic profile of these lexemes.

Table 3.24. Chi-squared (X^2) contributions for the related person/number features among the studied lexemes; statistically significant cells (with $df=3$) in boldface.

THINK.Z PERSON NUMBER\$cell.stat

Feature/Lexeme	<i>ajatella</i>	<i>miettiinä</i>	<i>pohtia</i>	<i>harkita</i>
Z_SG1	24.28	0.664	34.92	1.183
Z_SG2	1.940	19.51	29.66	7.665
Z_SG3	24.51	0.143	48.67	2.585
Z_PL1	1.711	0.352	4.348	1.868
Z_PL2	2.047	1.191	1.868	0.559
Z_PL3	2.463	10.49	0.273	1.588

Table 3.25. Significance values of the chi-squared (χ^2) contributions for the related person/number features among the studied lexemes; statistically significant cells (with $df=3$) in boldface.

THINK.Z_PERSON_NUMBER\$cell.p

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
Z_SG1	0.0605	1.000	- 2.52e⁻⁰³	1.000
Z_SG2	1.000	0.192	- 1.32e⁻⁰²	0.936
Z_SG3	0.0570	1.000	+ 1.99e⁻⁰⁵	1.000
Z_PL1	1.000	1.000	9.96e ⁻⁰¹	1.000
Z_PL2	1.000	1.000	1.000e ⁻⁰¹	1.000
Z_PL3	1.000	0.788	1.000	1.000

Table 3.26. Standardized Pearson residuals for the related person/number features among the studied lexemes; significant cells in boldface, i.e., $|e_{Feature, Lexeme}| > 2$.

THINK.Z_PERSON_NUMBER\$residual.pearson.std

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
Z_SG1	+7.64	-1.06	-7.48	-1.27
Z_SG2	+2.07	+5.54	-6.62	-3.10
Z_SG3	-9.07	-0.58	+10.4	+2.22
Z_PL1	+1.81	-0.69	-2.36	+1.43
Z_PL2	-2.01	+1.29	+1.57	-0.79
Z_PL3	+2.33	-4.05	+0.63	+1.41

Table 3.27. Simplified representation of the various methods of assessing cell-wise contributions for the distribution of the person/number features among the studied lexemes, with (+) denoting a significant observed deviation above the expected value, (-) a significant observed deviation below the expected value, and (0) a nonsignificant observed deviation.

THINK.Z_PERSON_NUMBER\$cell.sig

THINK.Z_PERSON_NUMBER\$residual.pearson.std.sig

Assessment strategy	Minimum significant value	Feature	ajatella	miettiä	pohtia	harkita
Table minimum	$\chi^2(df=15, \alpha=0.05) > 24.00$	Z_SG1	0	0	-	0
		Z_SG2	0	0	-	0
		Z_SG3	0	0	+	0
		Z_PL1	0	0	0	0
		Z_PL2	0	0	0	0
		Z_PL3	0	0	0	0
Cell-wise minimum	$\chi^2(df=1, \alpha=0.05) > 3.841$	Z_SG1	+	0	-	0
		Z_SG2	0	+	-	-
		Z_SG3	-	0	+	0
		Z_PL1	0	0	+	0
		Z_PL2	0	0	0	0
		Z_PL3	+	-	+	0
Standardized Pearson residual	$ e_{ij}/\text{standardized Pearson residual} > 2$	Z_SG1	+	0	-	0
		Z_SG2	+	+	-	-
		Z_SG3	-	0	+	+
		Z_PL1	0	0	-	0
		Z_PL2	-	0	0	0
		Z_PL3	+	-	0	0

The appropriate summary measures of association for the relationship between these six person/number features and the four lexemes are presented in Table 3.28. This time, Cramér's V at roughly 0.25 indicates that overall the person/number features are

not insignificant in accounting for the distribution of the studied lexemes. Furthermore, the feature-wise asymmetric association measures treating the lexemes as predictable dependents are only slightly higher than the opposite-direction associations; in addition, all of these measures, except both directions of λ are significant. Accordingly, in terms of interpretation, knowing the feature can be understood to allow us to account accurately for just below one-tenth of the behavior of the studied lexemes (as $\tau_{Lexeme|Feature} \approx 7.9\%$ and $U_{Lexeme|Feature} \approx 9.3\%$), whereas knowing the lexeme increases our accuracy in determining the feature by approximately 7% (as $\tau_{Feature|Lexeme} \approx 6.8\%$ and $U_{Feature|Lexeme} \approx 7.7\%$). Indeed, both of these measures of association, τ and U , whether calculated feature-wise or lexeme-wise for the person/number features, are quite small considering the possible range of $\hat{E}=[0,1]$. What is more, the association values were not much higher for the two major semantic classifications of AGENT presented earlier. So, at least in light of these two group-wise analyses, it would seem that association measures can be quite low at the same time as the observed distribution may be very significant, though naturally I cannot yet make a conclusive statement on this subject solely on the basis of these few example cases.

Table 3.28. The statistics (\hat{E}) and significance values (P-values) of the selected nominal measures of association, calculated both lexeme-wise, i.e., $\hat{E}(Feature|Lexeme)$, and feature-wise, i.e., $\hat{E}(Lexeme|Feature)$, for relationship between the six related person/number features and the four studied lexemes.

THINK.Z PERSON NUMBER\$associations

Measure	$\hat{E}(Feature Lexeme)$	$P(\hat{E})$	$\hat{E}(Lexeme Feature)$	$P(\hat{E})$
V	0.254	$2.16e^{-39}$	=	=
$\lambda_{A B}$	0.0107	0.7000	0.0227	0.442
$\tau_{A B}$	0.0684	$1.965e^{-75}$	0.0786	$1.20e^{-49}$
$U_{A B}$	0.0770	$5.580e^{-47}$	0.0929	$5.58e^{-47}$

As a final example of grouped analysis of closely related features we can take the semantic and structural classifications of another syntactic argument besides the AGENT. On the basis of the earlier descriptions of these lexemes, the syntactic PATIENT has been identified as the other major syntactic argument type of the studied lexemes in addition to the AGENT, and therefore its study is theoretically motivated and a useful supplement to the analyses of AGENT types among the studied lexemes. In fact, there are quite many more different types of PATIENTS than was the case with AGENTS. Not only do these include a large range of different semantic classifications of nominals (i.e., nouns and pronouns) as PATIENT arguments, but they also include different types of syntactic phrases and clauses, which is evident from the frequencies presented in Table 3.29. An analysis of the distribution presented in the simplified form presented above is given in Table 3.30.

Table 3.29. Contingency table presenting the frequencies of the occurrences of the different semantic and structural types of syntactic agents among the studied lexemes.
`multiple.feature.distribution(THINK.data, think.lex, SX_PAT.classes) %$ctab.ordered`

Feature/Lexeme	<i>ajatella</i>	<i>mieltiä</i>	<i>pohtia</i>	<i>harkita</i>	$\Sigma(\text{Feature} \text{Lexemes})$
SX_PAT.SEM_INDIVIDUAL	65	16	5	7	93
SX_PAT.SEM_GROUP	27	3	1	0	31
SX_PAT.SEM_NOTION	138	159	217	44	558
SX_PAT.SEM_ATTRIBUTE	18	18	26	5	67
SX_PAT.SEM_STATE	16	6	8	6	36
SX_PAT.SEM_TIME	21	7	8	2	38
SX_PAT.SEM_ACTIVITY	83	72	121	213	489
SX_PAT.SEM_EVENT	20	4	4	1	29
SX_PAT.SEM_COMM...	6	19	10	7	42
SX_PAT.SEM_COGNITION	8	6	2	2	18
SX_PAT.SEM_LOCATION	13	3	2	0	18
SX_PAT.SEM_ARTIFACT	12	1	1	2	16
SX_PAT.INDIRECT Q...	38	242	132	26	438
SX_PAT.DIRECT QUOTE	3	45	72	0	120
SX_PAT.INFINITIVE	38	0	1	3	42
SX_PAT.PARTICIPLE	65	0	3	6	74
SX_LX <i>että</i> CS.SX_PAT	317	48	23	8	396
$\Sigma(\text{Lexeme} \text{Features})$	888	649	636	332	2505

As we can see in Table 3.30, viewed lexeme-wise, *ajatella* is positively associated with INDIVIDUALS, GROUPS, TIME, EVENTS, and LOCATIONS, as well as INFINITIVES, PARTICIPLES, and *että*-clauses (equivalent to the English subordinate *that*-clauses) as syntactic PATIENTS. In contrast, *ajatella* is negatively associated with abstract NOTIONS, ACTIVITIES, and elements of COMMUNICATION, as well as INDIRECT QUESTIONS and DIRECT QUOTES as PATIENTS. For its part, *mieltiä* is positively associated with elements of COMMUNICATION as PATIENTS, in addition to INDIRECT QUESTIONS and DIRECT QUOTES, while it is negatively associated with GROUPS, INFINITIVES, PARTICIPLES, and *että*-clauses. Furthermore, *pohtia* is positively associated with abstract NOTIONS and ATTRIBUTES as well as INDIRECT QUESTIONS and DIRECT QUOTES as syntactic PATIENTS, whereas it is negatively associated with human INDIVIDUALS and GROUPS as well as INFINITIVES, PARTICIPLES, and *että*-clauses. Finally, *harkita* is positively associated with ACTIVITIES as PATIENTS, but negatively with human GROUPS and abstract NOTIONS, in addition to INDIRECT QUESTIONS, DIRECT QUOTES, and *että*-clauses. Taking the feature-wise angle, we can see that the strongest differentiating associations are those of ACTIVITIES with *harkita*, as well as INDIVIDUALS, GROUPS, and *että*-clauses with *ajatella*, in contrast with the other lexemes.

Table 3.30. Simplified representation of the various methods of assessing cell-wise contributions for the distribution of the different semantic and structural types of syntactic PATIENT arguments among the studied lexemes, with (+) denoting a significant observed deviation above the expected value, (-) a significant observed deviation below the expected value, and (0) a nonsignificant observed deviation.

multiple.feature.distribution(THINK.data,think.lex,SX_PAT.classes)\$residual.pearson.std.sig

Feature/Lexeme	ajatella	miettä	pohtia	harkita
SX PAT.SEM INDIVIDUAL	+	0	-	0
SX PAT.SEM GROUP	+	-	-	-
SX PAT.SEM NOTION	-	0	+	-
SX PAT.SEM ATTRIBUTE	0	0	+	0
SX PAT.SEM STATE	0	0	0	0
SX PAT.SEM TIME	+	0	0	0
SX PAT.SEM ACTIVITY	-	-	0	+
SX PAT.SEM EVENT	+	0	0	0
SX PAT.SEM COMMUNICATION	-	+	0	0
SX PAT.SEM COGNITION	0	0	0	0
SX PAT.SEM LOCATION	+	0	0	0
SX PAT.SEM ARTIFACT	+	0	0	0
SX PAT.INDIRECT QUESTION	-	+	+	-
SX PAT.DIRECT QUOTE	-	+	+	-
SX PAT.INFINITIVE	+	-	-	0
SX PAT.PARTICIPLE	+	-	-	0
SX LX että CS.SX PAT	+	-	-	-

The appropriate summary measures of association for the relationship between these different types of syntactic agents and the four lexemes are presented in Table 3.31. This time, Cramér's V , at as high as roughly 0.45, is a clear indication that overall the types of PATIENTS have a very important role in the use of the studied lexemes, and they can be seen to account for the distribution of the studied lexemes; the *Effect Size* is even higher as $w=0.775$. In contrast to the other two example cases, all of the three different association measures, including the λ , are significant. Furthermore, all the feature-wise asymmetric association measures, treating the lexemes as predictable dependents, are approximately twice as high as the respective values for the opposite-direction associations. Thus, knowing the semantic or structural classification of the PATIENT accounts for one-fifth of variation of the lexeme, with $\tau_{Lexeme|Feature}=0.215$ and $U_{Lexeme|Feature}=0.216$. In contrast, knowing the lexeme explains roughly one-tenth of the different feature types of PATIENTS, with $\tau_{Feature|Lexeme}=0.0978$ and $U_{Feature|Lexeme}=0.131$. So, in the case of the different types of PATIENTS, the feature-wise association measure values are substantially higher than the lexeme-wise ones, which is contrary to what was the case with the AGENTS and person/number features.

Table 3.31. The statistics (\hat{E}) and significance values (P-values) of the selected nominal measures of association, calculated both lexeme-wise, i.e., $\hat{E}(Feature|Lexeme)$, and feature-wise, i.e., $\hat{E}(Lexeme|Feature)$, for relationship between the different types of syntactic PATIENTS and the four studied lexemes.

THINK.SX PAT.SEM ALL\$associations

Measure	$\hat{E}(Feature Lexeme)$	$P(\hat{E})$	$\hat{E}(Lexeme Feature)$	$P(\hat{E})$
V	0.448	$6.40e^{-284}$	=	=
$\lambda_{A B}$	0.221	$1.54e^{-48}$	0.311	$2.82e^{-70}$
$\tau_{A B}$	0.0978	0.0	0.215	$1.42e^{-306}$
$U_{A B}$	0.131	$1.60e^{-270}$	0.216	$1.60e^{-270}$

These examples of the semantic and other classifications of the AGENTS and PATIENTS as well as the person/number features of the studied lexemes have shown that much insight can be gained by the grouped study of closely related features in the manner shown above. However, one should note that this type of scrutiny does not consider the relationships and interactions of a set of related features with other individual features or their sets which may also occur in the context of the studied lexemes. Therefore, this set-wise analysis does not do away with the need for full-scale multivariate methods, though it is quite informative in itself. Furthermore, the singular feature analyses are still useful and necessary in selecting those individual features which are substantial and significant to the degree that they should be included in the scrutiny with full-scale multivariate methods. But before reaching that stage, it is first worthwhile (and necessary) to observe and scrutinize their pairwise co-occurrences and interactions.

3.3 Bivariate methods

3.3.1 General considerations

Until this point I have focused on the relationship of individual contextual features, or to a lesser extent, sets of closely related and complementary features, with the studied lexemes. However, a large proportion of the selected features can at least theoretically co-occur with each other. That is to say, there is nothing in the structure of the linguistic analysis and description scheme that I follow that inherently blocks their joint occurrence, though in practice some of these feature combinations may be rare or non-occurrent due to semantic, pragmatic or other considerations which our present descriptive apparatus does not yet fully account for. It is therefore of linguistic interest to scrutinize pairwise the selected features, in order to observe the degree to which they jointly occur, or do not, among the studied lexemes.

This pairwise analysis will indicate self-evident associations due to overlap, explicitness, and redundancy in our descriptive system. These are due to 1) logical (symmetric) complementarity of the type studied above in Section 3.2.3, such as all verbs being either FINITE or NON-FINITE but not both at the same time, 2) directional compositionality, i.e., all infinitives are NON-FINITE (but not all NON-FINITE forms are participles), or simply 3) overlap, i.e., a FINITE verb with an overt subject/AGENT must *per definitionem* be in an ACTIVE form. However, pairwise scrutiny can also reveal non-obvious linguistic preferences and potentially idiomatic constructions. Furthermore, this stage is useful, and, in fact, necessary in identifying those features that correlate with each other to the extent that it has to be taken into consideration in the successful application of the subsequent multivariate methods.

To make things simpler, the pairwise comparisons will be based on the methods already presented in the singular feature analyses above. Here, however, the two variables under scrutiny are not an individual feature (or related set of features) on the one hand and the set of studied lexemes on the other, but two distinct features instead, which are assessed in the simplest case in terms of their joint or partial occurrences or nonoccurrences in the data. In this setting, the perfect positive pairwise relationship would firstly mean that the occurrence of one feature is always matched by the occurrence of the other feature, both ways, and secondly that the nonoccurrence of either feature is always matched by the nonoccurrence of the other feature. In contrast, a perfect negative pairwise relationship would entail that the occurrence of one feature would always imply the nonoccurrence of the other, and vice versa. However, these require that the frequencies of both features are equal, which we know not to be the case for the most part. Nevertheless, we are interested in evaluating both the strength of the overall relationship between any two features, and furthermore, the strength of the directional relationships. In other words, does knowing the occurrence or nonoccurrence of one feature allow us to determine the occurrence or nonoccurrence of the other feature, and to what extent this is the case. These are exactly the types of questions that we can address with summary measures of association, already presented above in Section 3.2.2 among the univariate methods.

3.3.2 Pairwise associations of individual features

Let us take as an example two of the features that we have already studied individually, namely, the FIRST PERSON SINGULAR (*Z_SG1*) as a morphological feature of the studied lexemes and the human INDIVIDUAL as a semantic type of their syntactic AGENTS (*SX_AGE.SEM_INDIVIDUAL*). From the outset we may suspect that there should be substantial overlap, which we can systematically assess with the help of Table 3.32. We can see that the two features in question co-occur 246 times, and furthermore, that the *Z_SG1* feature (almost⁵⁸) always occurs with an INDIVIDUAL AGENT, as can be logically expected. However, not all INDIVIDUAL AGENTS are FIRST PERSON SINGULAR forms (represented by 2005 instances), at least in the data we use. This is not really surprising as all of the six different person/number features, of which the FIRST PERSON SINGULAR is but one, are by definition classified as INDIVIDUAL AGENTS, regardless of whether they have an overt AGENT or not. Furthermore, there is a total of 1151 instances in the data with neither of the two features in question occurring in the context of the studied lexemes. The summary measures of association representing the pairwise relationship between these two features are presented in Table 3.33.

Table 3.32. The joint distribution of the *SX_AGE.SEM_INDIVIDUAL* feature and the *Z_SG1* feature among the studied lexemes.

```
singular.pairwise.association(cbind(THINK.data["SX_AGE.SEM_INDIVIDUAL"], THINK.data["Z_SG1"]))
```

Feature ₁ /Feature ₂	<i>Z_SG1</i>	<i>-Z_SG1</i>	Σ(Row)
<i>SX_AGE.SEM_INDIVIDUAL</i>	246	2005	2251
<i>¬SX_AGE.SEM_INDIVIDUAL</i>	2	1151	1153
Σ(Column)	248	3156	3404

For the assessment of the overall pairwise relationship we can use Cramér's V , which is 0.195 for these two features. Furthermore, this value is very significant, implying a real relationship between the two features. For the directional assessment of the pairwise relationship we can, in principle, use any of the asymmetric measures. Of these, the earlier in-depth comparisons of the various available methods would indicate that both the Goodman-Kruskal $\tau_{A|B}$ and Theil's Uncertainty Coefficient $U_{A|B}$ would be the best ones, with a slight preference for the latter of the two. However, for 2x2 tables as is the case here, the value of $\tau_{A|B}$ is by definition the same in both directions, that is, it becomes (only) in such a particular setting a symmetric measure (see Costner 1965: 351), whereas Theil's $U_{A|B}$ retains its asymmetry. Therefore, the Uncertainty Coefficient $U_{A|B}$ becomes slightly more preferable, as the potential differences of its two asymmetric versions allow us to evaluate the directionality of the pairwise relationship. As we can see, knowing that a studied lexeme has (or does

⁵⁸ The two non-INDIVIDUAL cases of the *Z_SG1* feature are in fact errors that remained in the data even after the automatically parsed analysis had repeatedly been combed through manually. These errors were discovered only at this late stage of reporting the results. Specifically, the underlying form in question is *ajatellen*, which can be morphologically analyzed as either the INSTRUCTIVE case of the SECOND INFINITIVE or the FIRST PERSON SINGULAR of the POTENTIAL mood of *ajatella*. Of the two, the former analysis is correct for these two cases, and probably in general, too. Of course, I could have corrected these two cases, but I chose instead to leave them as a demonstration of the possible sources of error in linguistic data analysis, and furthermore as an example that such occasional errors will not have a significant bearing on the overall analysis, when the sample is sufficiently large as is the case here.

not have) an INDIVIDUAL AGENT allows us to determine whether the studied lexeme is (or is not) in the FIRST PERSON SINGULAR form with $U_{Z_SG1|SX_AGE.SEM_INDIVIDUAL}=0.109$. This is more than twice as much as in the opposite direction, with $U_{SX_AGE.SEM_INDIVIDUAL|Z_SG1}=0.0445$. This is in accordance with the logical directionality of the FIRST PERSON SINGULAR feature being subsumed by the INDIVIDUAL type of AGENT. That this particular pairwise relationship at best accounts for only about 10% of the overall variation of the studied lexemes is, in addition, due to the fact that roughly one-third (33.8%) of the studied lexemes do not occur with either of the two contextual features.

Table 3.33. Values of selected measures of association for the evaluation of the pairwise relationship between the SX_AGE.SEM_INDIVIDUAL and the Z_SG1 features among the studied lexemes.

Association measure ($\hat{E}_{Feature\ 1 Feature\ 2}$)	Value	Significance (P-value)
Cramér's V	0.195	$6.86e^{-30}$
$\tau_{Z_SG1 SX_AGE.SEM_INDIVIDUAL}$	0.0384	$3.14e^{-30}$
$\tau_{SX_AGE.SEM_INDIVIDUAL Z_SG1}$	0.0384	$3.14e^{-30}$
$U_{Z_SG1 SX_AGE.SEM_INDIVIDUAL}$	0.109	$4.38e^{-44}$
$U_{SX_AGE.SEM_INDIVIDUAL Z_SG1}$	0.0445	$4.38e^{-44}$

As an example of the pairwise comparison of logically complementary features we can take the two already studied semantic types of AGENTS, namely, human INDIVIDUALS and human GROUPS, denoted by the labels SX_AGE.SEM_GROUP and (SX_AGE.SEM_INDIVIDUAL, respectively. This is somewhat artificial as an example, since we know from the outset that their distribution is complementary, which can also be clearly seen in the joint distribution of their occurrences and nonoccurrences presented in Table 3.34 and, to a lesser extent, in the summary measures of association in Table 3.35. There are no common occurrences, as should naturally be the case since an AGENT in the classification scheme used in this study can have only one semantic classification. Furthermore, the overall relationship between the two features has a relatively high value of Cramér's V at 0.397, which is significant without a doubt. Accordingly, knowing that a studied lexeme has an INDIVIDUAL as its AGENT allows us to determine that the AGENT cannot be a GROUP, with $U_{SX_AGE.SEM_GROUP|SX_AGE.SEM_INDIVIDUAL}=0.328$, which is more than twice the corresponding value in the opposite direction, i.e., $U_{SX_AGE.SEM_INDIVIDUAL|X_AGE.SEM_GROUP}=0.137$. Again, this clearly complementary but less than perfect negative relationship is explained by the substantial number (897 instances, i.e., 26.4%) of studied lexemes without either semantic type of AGENT, implying that these lexemes have no AGENT at all. Knowing the syntactic and morphological general characteristics of Finnish verbs, I can make an educated guess that these cases are most probably forms in the PASSIVE voice or NON-FINITE PARTICIPIAL or INFINITIVAL (CLAUSE-EQUIVALENT) forms.

Table 3.34. The joint distribution of the SX_AGE.SEM_INDIVIDUAL feature and the SX_AGE.SEM_GROUP feature among the studied lexemes.
`singular.pairwise.association(cbind(THINK.data["SX_AGE.SEM_INDIVIDUAL"], THINK.data["SX_AGE.SEM_GROUP"]))`

Feature ₁ /Feature ₂	SX_AGE.SEM_GROUP	¬SX_AGE.SEM_GROUP	Σ(Row)
SX_AGE.SEM_INDIVIDUAL	0	2251	2251
¬SX_AGE.SEM_INDIVIDUAL	256	897	1153
Σ(Column)	256	3148	3404

Table 3.35. Values of selected measures of association for the evaluation of the pairwise relationship between the SX_AGE.SEM_INDIVIDUAL and the SX_AGE.SEM_GROUP features among the studied lexemes.
`singular.pairwise.association(cbind(THINK["SX_AGE.SEM_INDIVIDUAL"], THINK["SX_AGE.SEM_GROUP"]))`

Association measure ($\hat{E}_{Feature 1 Feature 2}$)	Value	Significance (P-value)
Cramér's V	0.397	$7.50e^{-119}$
$\tau_{SX_AGE.SEM_GROUP SX_AGE.SEM_INDIVIDUAL}$	0.159	$1.64e^{-119}$
$\tau_{SX_AGE.SEM_INDIVIDUAL SX_AGE.SEM_GROUP}$	0.159	$1.64e^{-119}$
$U_{SX_AGE.SEM_GROUP SX_AGE.SEM_INDIVIDUAL}$	0.328	$1.17e^{-131}$
$U_{SX_AGE.SEM_INDIVIDUAL SX_AGE.SEM_GROUP}$	0.137	$1.17e^{-131}$

We will get a better overview of the pairwise relationships when we scrutinize individual pairings in relation to all the rest, which will be presented in Section 4.2.1 with the bivariate results to follow below. Lacking a natural threshold in pruning excessively correlating features, I will nevertheless resort to the general ones presented above in Section 3.2.2. Thus, when the relationship for a feature pairing is by all accounts *strong*, that is, when the value of the association measure exceeds $U_{A|B} > 0.5$ at least in one direction, and therefore, at least one of the features accounts for a majority of the variation of the other, I will in such a case include only one of the two features in question into the multivariate analysis. Nevertheless, this task must be undertaken from an overall perspective with a linguistically informed, careful consideration of the entire feature set to be selected. In addition, I will also scrutinize pairings exhibiting a *moderate* relationship, i.e., $U_{A|B} > 0.2$, as such associations may also turn out to be of some interest. Moreover, the overall pairwise results will also allow us to evaluate the value range of mutual pairwise associations among the features to be included in the final multivariate analysis, thus giving us some idea of the level of multicollinearity among them.

3.3.3 Pairwise comparisons of two sets of related features

In addition to these pairwise comparisons, we could quite naturally be interested in the relationships and joint interaction of more than two features. In principle, this can be done, but for the sake of methodological simplicity, I will here limit the study to a bivariate analysis. However, we can make an extension of these pairwise comparisons of singular individual contextual features to the simultaneous study of two sets of closely related (complementary) features. These sets of features can be treated as different values (or, classes or categories) of the two general variables and analyzed in a manner very similar to what was done above in Section 3.2.3. For instance, we could be interested in the pairwise relationship between the different semantic types

of AGENTS and the PATIENTS, which I have already studied individually. So, the joint distributions of the semantic and structural types of syntactic AGENTS and PATIENTS are presented in Table 3.36, and the results of the ensuing analysis are shown in simplified form in Table 3.37. Only the very rarest semantic categories of PATIENTS have been left out, namely, SUBSTANCES (2 instances), FOOD (2), FLORA (1), the BODY (1), amounting to 6 instances in all (corresponding to only 0.2% of the altogether 2812 instances PATIENT arguments).⁵⁹

Table 3.36. Contingency table presenting the frequencies of the joint occurrences of the different semantic and structural types of syntactic AGENTS and PATIENTS among the studied lexemes.

THINK.SX AGE.SX PAT\$ctab.ordered

Patient/Agent (SX PAT/SX AGE)	SEM_INDIVIDUAL	SEM_GROUP	Σ(Patient)
SEM_INDIVIDUAL	65	5	70
SEM_GROUP	18	2	20
SEM_NOTION	316	60	376
SEM_ATTRIBUTE	39	3	42
SEM_STATE	17	3	20
SEM_TIME	20	4	24
SEM_ACTIVITY	225	90	315
SEM_EVENT	7	1	8
SEM_COMMUNICATION	30	1	31
SEM_COGNITION	12	0	12
SEM_LOCATION	6	1	7
SEM_ARTIFACT	10	0	10
INDIRECT_QUESTION	330	37	367
DIRECT_QUOTE	119	1	120
INFINITIVE	34	3	37
PARTICIPLE	53	5	58
SX_LX_että_CS	324	7	331
Σ(Agent)	1625	223	1848

⁵⁹ One could consider collapsing these and some of the other less frequent categories into the more frequent ones, e.g., ATTRIBUTE and STATE as subtypes of abstract NOTION. However, as all the semantic categories here belong to the top-level unique beginners in the WordNet ontology, one might in the resultant supersets lose in their internal coherence what one would benefit from the decrease in the number of variables.

Table 3.37. Simplified representation of the cell-wise contributions for the joint distribution of the different semantic and structural types of syntactic AGENT and PATIENT arguments among the studied lexemes using standardized Pearson residuals, with (+) denoting a significant observed deviation above the expected value, (-) a significant observed deviation below the expected value, and (0) a nonsignificant observed deviation.

THINK.SX AGE.SX PAT\$residual.pearson.std.sig

Patient/Agent (SX_PAT.SEM_XXX/ SX_AGE.SEM_XXX)	SX_AGE. SEM_INDIVIDUAL	SX.AGE. SEM_GROUP
SX_PAT.SEM_INDIVIDUAL	0	0
SX_PAT.SEM_GROUP	0	0
SX_PAT.SEM_NOTION	-	+
SX_PAT.SEM_ATTRIBUTE	0	0
SX_PAT.SEM_STATE	0	0
SX_PAT.SEM_TIME	0	0
SX_PAT.SEM_ACTIVITY	-	+
SX_PAT.SEM_EVENT	0	0
SX_PAT.SEM_COMMUNICATION	0	0
SX_PAT.SEM_COGNITION	0	0
SX_PAT.SEM_LOCATION	0	0
SX_PAT.SEM_ARTIFACT	0	0
SX_PAT.INDIRECT_QUESTION	0	0
SX_PAT.DIRECT_QUOTE	+	-
SX_PAT.INFINITIVE	0	0
SX_PAT.PARTICIPLE	0	0
SX_LX_että_CS.SX_PAT	+	-

We can see from Table 3.36 that the joint occurrences of the selected different semantic and structural types of AGENTS and PATIENTS (1848 instances) account for almost all (91.9%) of the joint occurrences of both argument types (2011 instances). However, these joint occurrences of both AGENT and PATIENT types constitute barely a majority (57.2%) of the individual overall frequencies of either argument type with the studied lexemes (altogether 3231 instances). Incidentally, this last figure also means that practically all (94.9% of the overall total 3404) of the studied lexemes have either an AGENT, a PATIENT, or both as an argument. Nevertheless, taking into account the overall marginal frequencies for each feature type in the cell-wise assessment, only a few of the AGENT/PATIENT type combinations exhibit a significant deviation. Thus, GROUP AGENTS and abstract NOTION or ACTIVITY PATIENTS as well as INDIVIDUAL AGENTS and DIRECT QUOTES or *että*-clauses are positively associated, whereas INDIVIDUAL AGENTS and abstract NOTION or ACTIVITY PATIENTS as well as GROUP AGENTS and DIRECT QUOTES or *että*-clauses are negatively associated with each other. Looking at the raw frequency data in Table 3.36, some of the AGENT/PATIENT combinations such as INDIVIDUAL and GROUP PATIENTS are clearly more frequent in absolute terms with INDIVIDUAL instead of GROUP AGENTS. However, in terms of proportions of such PATIENT types with respect to the two AGENTS types, the differences are, nonetheless, not significant (enough) to show up.

Table 3.38. The statistics (\hat{E}) and significance values (P-values) of the selected nominal measures of association, calculated both lexeme-wise, i.e., $\hat{E}(Feature|Lexeme)$, and feature-wise, i.e., $\hat{E}(Lexeme|Feature)$, for the pairwise relationship between the different types of syntactic AGENTS and PATIENTS among the four studied lexemes.

THINK.SX AGE.SX PAT\$associations

Association measure ($\hat{E}_{Feature\ 1 Feature\ 2}$)	Value	Significance (P-value)
Cramér's V	0.278	$2.26e^{-22}$
$\tau_{PATIENT AGENT}$	0.0139	$2.90e^{-77}$
$\tau_{AGENT PATIENT}$	0.0771	$2.34e^{-22}$
$U_{PATIENT AGENT}$	0.0188	$7.56e^{-24}$
$U_{AGENT PATIENT}$	0.110	$7.56e^{-24}$

In terms of summary measures of association, the overall association of the various types of the two arguments is substantial, with Cramér's V at 0.2778; this association is significant and equal to the *Effect Size* (as the underlying table is of the form 2xN). As the number of different types of AGENTS is substantially less than the number of PATIENT types, it is no surprise that the asymmetric measures with PATIENT as the independent dimension and AGENT as the dependent, predicted one are many times greater, whether measured in terms of $\tau_{AGENT|PATIENT}=0.0771$ or $U_{AGENT|PATIENT}=0.110$, than the opposite-direction measures $\tau_{PATIENT|AGENT}=0.0139$ or $U_{PATIENT|AGENT}=0.0188$. Nevertheless, all of these measures indicate that they account for at most one-tenth of the variation of the studied lexemes.

In conjunction with pairwise comparisons, Gries (2003a: 101-106) suggests assessing the strength of individual features against the rest when their preferences of association are in conflict. This concerns in Gries' dichotomous setting cases where two features are observed to co-occur, but the overall preferences of these two features differ, in that the first feature is positively associated with one form of the construction whereas the other feature is positively associated with the alternative construction. Gries (2003a: 130, note 25) proposes counting the occurrences of a feature with a positive association with one form of the construction against all its co-occurrences with features having a negative association, and then calculating the overall index ratio to discover which of the two alternative constructions prevails.

This type of analysis can be extended to the polytomous case of the four alternative lexemes studied here by counting lexeme-wise for each feature positively associated with that lexeme the co-occurrences of this particular feature with all the features negatively associated with the same lexeme. This results in a 2x2 contingency table with a generic structure presented in Table 3.39. Then, for each positively associated feature we can calculate the ratio of occurrences of the preferred lexeme against the other lexemes. Furthermore, we can evaluate the overall strength of each relationship represented in such a contingency table with the symmetric Cramér's V , and we can also calculate a significance level for this comparison. With four possible alternative lexemes, however, it is quite probable that the ratios will generally be smaller (and even negative), in comparison to Gries' (2003a) setting with only two alternatives. In addition to contrasting individual positively associated features against negatively associated ones, we can just as well calculate, in a similar fashion, the opposite case of individual negatively associated features against the positively associated ones, as well as the relative weights of individual features against all the rest among the positively associated features, or individual features against all the rest among all the

negatively associated ones. Nevertheless, for reasons of space I will not pursue this analysis strategy further in this dissertation.

Table 3.39. Table representing the lexeme-wise adaptation of Gries' (2003a) proposal for calculating the relative weight of an individual positively associated feature against all the co-occurring, overall negatively associated features, where $F_{positive|Lexeme}(i)$ is an individual positively associated feature for some lexeme, $F_{positive|Lexeme}$ is the entire set of features positively associated with some lexeme, and $F_{negative|Lexeme}$ is the entire set of positively associated features for the same lexeme.

Joint conditions: Positive against negative features	Lexeme	¬Lexeme
$F_{positive Lexeme}(i) \wedge F_{negative Lexeme}$	$\sum O[F_{positive Lexeme}(i) \wedge F_{negative Lexeme} \wedge Lexeme]$	$\sum O[F_{positive Lexeme}(i) \wedge F_{negative Lexeme} \wedge \neg Lexeme]$
$\neg F_{positive Lexeme} \wedge F_{negative Lexeme}$	$\sum O[\neg F_{positive Lexeme}(i) \wedge F_{negative Lexeme} \wedge Lexeme]$	$\sum O[\neg F_{positive Lexeme}(i) \wedge F_{negative Lexeme} \wedge \neg Lexeme]$

3.4 Multivariate methods

3.4.1 Logistic regression analysis with nominal outcomes and variables

The general purpose of multivariate analysis is to study the joint and simultaneous relationship of all the selected variables with respect to the studied phenomenon. In this linguistic study, the key question is the relationship of the contextual features with the studied four lexemes (which are all nominal variables). Though there are many possible foci of interest in multivariate analysis, I am primarily concerned with only two. Firstly, I am interested in the relative weights and differences in the impact of the individual variables which have been identified as pertinent in the preceding univariate and bivariate analyses. Secondly, I also wish to know how well the selected variables are able to explain and account overall for the linguistic phenomenon under study. This relationship between the lexemes and features can be considered directionally skewed, since for each observed instance in the data only one of the four lexemes at a time is associated with a varying (but potentially greater) number of features present in the context. In this setting, it makes more sense to study which one of the studied lexemes can be expected to occur, given a particular context constituted by some set of features, than the other way around. Furthermore, from prior research (e.g., Arppe and Järvikivi 2007b, Featherston 2005) and from the univariate examples in Section 3.2, we know that in practice individual features or sets of features are *not* observed in corpora to be categorically matched with the occurrence of only one of the lexemes in a synonymous set and none of the others. Rather, while one lexeme in a synonymous set may be by far the most frequent in some particular context, others also do occur, albeit with often a considerably lower relative frequency.

In addition, these earlier studies indicate that even though the observed relative frequency differences may be very great, in acceptability ratings by native speakers some of the less frequent alternatives can be almost as highly rated as the most frequent one. In other words, in terms of acceptability alternative linguistic items for semantically similar content, whether syntactic constructions or lexemes in some synonym set, are arranged along a gradual continuum instead of a dichotomy. With this in mind, the representation of linguistic reality in multivariate analysis is probably more accurate when we reformulate the relationship between lexemes and contextual features, so that we rather study the *expected probabilities of occurrence* of all the individual lexemes belonging to a synonymous set, given some contextual features, instead of a discrete choice of only one of the four alternative lexemes, allowing only for the dichotomous values of occurrence or nonoccurrence. That is, we in effect shift our focus from individual instances of discrete choices in usage to the overall observations in the data, where expected probability can be understood in terms of the proportions of occurrences of each lexeme of the synonymous set, given a set of contextual features.

For this purpose, *multinomial* (alternatively also referred to as *multiple-category*, *multiple-class*, *polytomous*, *polychotomous*, or even, *discrete-choice*) *logistic regression* analysis (see Fox 1997: 467-472; Hosmer and Lemeshow 2000: 260-287; Agresti 2002: 267-274; Cohen et al. 2003: 519-522; Menard 1995: 80-86) is an attractive multivariate statistical approach. A proper multinomial logistic regression model for K outcomes is based on the simultaneous, joint fitting of a set of $K-1$ simpler *binary logistic regression* models (originating from Cox 1958; see also

Harrell 2007: 215-267) against some baseline category (denoted here as class K). If no algorithm is available for implementing such joint fitting, or if designating a baseline category turns out to be problematic for other reasons, there are various heuristics that approximate it through differing partitions of the multi-class model into sets of binary logistic regression models which can then be fitted individually, independently of each other. The general scheme for representing the $K-1$ formulas of a multinomial logistic model for K classes for outcome Y , with class K set as the baseline, resulting from the joint effect of a set X of M explanatory variables selected in the model and having a set β of $(K-1) \cdot M$ parameters and a set α of $(K-1)$ constant intercepts, is presented in formulas 3.17-3.19 below. Specifically, the selection of M explanatory variables $X = \{X_1, \dots, X_M\}$ is then understood to constitute the *model* of the studied phenomenon. Sometimes, the explanatory variables are alternatively referred to as *predictors*, and the associated parameters as *coefficients* of the model.

$$(3.17) P_k(X) = P(Y=k|X), \text{ with } \sum_{k=1 \dots K} P_k(X) = 1 \text{ and } k = \{1, \dots, K\}, \text{ and } P_K(X) = P(Y=K|X) = 1 - \sum_{k=1 \dots K-1} P_k(X) \text{ as the baseline case.}$$

$$(3.18) \log_e[P_k(X)/P_K(X)] = \alpha_k + \beta_k X \Leftrightarrow P_k(X) = \exp(\alpha_k + \beta_k X) / [1 + \sum_{k=1 \dots K-1} \exp(\alpha_k + \beta_k X)] \text{ for } k=1 \dots K-1 \text{ and } P_K(X) = 1 - \sum_{k=1 \dots K-1} P_k(X) \text{ (the baseline thus assigned the "leftover" probability)}$$

$$(3.19) \beta_k X = \beta_{k,1} X_1 + \beta_{k,2} X_2 + \dots + \beta_{k,M} X_M$$

with classes $k = \{1, \dots, K-1\}$, and M explanatory variables $X = \{X_1, \dots, X_M\}$, parameters $\beta = \{(\beta_{1,1}, \dots, \beta_{1,M}), (\beta_{2,1}, \dots, \beta_{2,M}), \dots, (\beta_{K-1,1}, \dots, \beta_{K-1,M})\}$, and constants $\alpha = \{\alpha_1, \dots, \alpha_{K-1}\}$.

As a *direct probability model* (Harrell 2001: 217), multinomial as well as binary regression yields probability estimates, corresponding to the expected proportions of occurrences, conditional on the values of the explanatory variables that have been selected for inclusion in the models. Most crucially, multinomial logistic regression and its various approximations (via the binary logistic functions they are based on) are applicable for nominal variables such as the contextual features used in this study, for which it has a natural interpretation concerning their effect on the outcome probabilities. In general, for any type of variable included in the multinomial model, once the parameters β have been fitted with the data, each specific parameter (fitted *coefficient*) $\beta_{k,m}$, associated with each variable X_m in each of the constituent non-baseline binary models, can be interpreted as the *logarithm of the odds* (known also as *log-odds* or *logits*) per unit change of the particular variable that the outcome is a given particular class k , in comparison to the baseline class K – with the other variables being equal and with no interactions assumed. The actual odds are then equal to the base of the natural logarithm e to the power of $\beta_{k,m}$, i.e., $e^{\beta_{k,m}}$ (see Harrell 2001: 218, equation 10.11). These odds can also be formulated as the ratio of the probabilities of class k occurring in comparison to class K occurring, and these probabilities can again be understood as proportions of overall occurrences in the data.

As has already been done starting with the univariate analysis above, the explanatory variables to be included in the multinomial logistic regression analysis in this study are individual contextual features, which have the logical value TRUE when occurrent in the context in a relevant way (i.e., belonging to the syntactic argument structure or

morphological make-up of the instances of the studied lexemes), and the value FALSE when this is not the case. In fact, in some situations the value can be FALSE even if the feature in question is not applicable at all in the particular context and could thus not occur even in principle. For computational purposes, these two logical values can be represented as 1 and 0, respectively.⁶⁰ In terms of interpretation, when the feature represented by the variable is present in the context and the associated variable thus switched to TRUE instead of FALSE, the parameter $\beta_{k,m}$ for each such binary nominal explanatory variable X_m is the associated increase in the logarithm of the odds (i.e., log-odds) of the outcome belonging to a selected class k in comparison with the baseline category K , with the other explanatory variables remaining equal.

What this means in practice is that, if the parameter (coefficient) and thus the log-odds for some hypothetical class k and nominal binary variable X_m is, e.g., $\beta_{k,m}=2$, the odds of the outcome being class k in comparison to the baseline class K is $e^2 \approx 7.4 \sim 37:5$ when the associated variable is TRUE. In other words, it would be over seven times more likely to encounter class k than the baseline class K , when the feature is to be found in the context, other things being equal. At the same time, however, we could also expect the baseline class K to occur with the inverse ratio of $1/e^2 \approx 0.14 \sim 5:37$, that is, approximately once in every eight times that the feature is present in the context. If the parameter is $\beta_{k,m}=0$, the explanatory variable in question would not have a substantial bearing on the outcomes (in comparison to the other variables), since the odds would be $e^0=1 \sim 1:1$. In contrast, if the parameter were negative such as $\beta_{k,m}=-1$, the odds in such a case would be against class k and thus, in favor of the baseline class K , with $e^{-1}=0.38 \sim 3:8$ (see Harrell 2001: 217-220; Agresti 2002: 166-167; Cohen et al. 2003: 492-493). However, one should note that the *nonoccurrence* of a feature in the context, with the associated explanatory variable thus being FALSE, does *not* in itself give us any information about the odds of any of the outcomes, that is, lexemes occurring. That is, the odds apply *only* when the feature in question is actually present.

Furthermore, for each parameter $\beta_{k,m}$ its asymptotic standard-error (*ASE*) can be calculated,⁶¹ which can then be used to assess the significance of the parameter in question in comparison the null hypothesis, that is, that the particular variable had no effect and that the parameter would thus equal zero (according to formulas 3.20-3.21 below, see Fox 1997: 450). Alternatively, the *ASE* can be used to calculate a Confidence Interval *CI* for the parameter, which is then significant if the confidence interval does not include zero (formula 3.22 below, Cohen et al. 2003: 497-498).

$$(3.20) z = \beta_{k,m} / ASE$$

$$(3.21) P(\beta_{k,m} \neq 0) = P(|Z| > |z|) = 2 \cdot [P_{N(0,1)}(|z|)]$$

$$(3.22) \text{Confidence Interval } CI = \beta_{k,m} \pm z_{1-\alpha/2} \cdot ASE, \text{ when } z \sim N(0,1)$$

⁶⁰ This is in accordance with *dummy-variable coding*, which is the most common coding convention. However, there are also other types of possible coding schemes (see, e.g., Cohen et al. 2003: 302-253).

⁶¹ Fortunately, these are calculated by statistical software as part of the fitting process of the logistic regression model. For instance, in the *R* statistical environment the `glm` function automatically calculates not only the *ASE* but also the standardized coefficient (z) and the associated P-level. These can be specifically accessed via the function call sequence `coef(summary(glm(...)))`.

In similar linguistic settings, I have not encountered the use of multinomial logistic regression, or its approximations, other than my own exploratory study (Arppe 2006c) and the preliminary reporting of the results of this dissertation in Arppe (2007). However, the simpler basic method for binary logistic regression has been used by, for example, Bresnan et al. (2007) in the study of the English dative alternation, and by Grondelaers et al. (2002) in the study of the occurrence of the Dutch *er* ‘there’.

3.4.2 Selection of variables in multivariate logistic regression

As was noted above, multinomial and binary logistic regression analysis is based on constructing a model which consists of the individual explanatory variables and their interactions, which are hypothesized to explain the studied phenomenon and determine the probabilities of the associated outcomes, whether with the observed data used to fit the model or future data to be predicted by the model. In practice, the maximum number of variables (including their interaction terms) that can produce a valid and reliable model is limited by and proportionate to the size of the available data. If there are too many variables in relation to the data, the resultant model will increasingly represent noise and spurious relations rather than real effects between the outcomes and the explanatory variables. This is called *overfitting* the model, meaning that the model will fit the data at hand and its idiosyncrasies *too* well and will consequently generalize poorly to unseen new data. Rules of the thumb have been presented for *limiting sample sizes* (m), with which the maximum recommended number of variables p_{max} is proportionate, being approximately between $m/10$ and $m/20$.

For binary logistic regression models, of which the multinomial models studied here consist, the limiting sample size $m = \min(n_1, n_2)$, where n_1 and n_2 are the overall frequencies of the two alternative outcomes (Peduzzi et al. 1996; see also Hosmer and Lemeshow 2000: 346-347; Harrell 2001: 60-61, Table 4.1). In this study with four alternative synonyms, the minimum outcome frequency of any of the possible component binary logistic regression models is equal to the overall frequency of the least frequent of the studied four lexemes, that is, *harkita* with 387 occurrences, which thus becomes also the limiting sample size. Therefore, in order to avoid overfitting, the number of explanatory variables to be included in the multivariate model would be restricted to approximately $387/10 \approx 39$, say around 40 at the most. This may appear a conservative limitation as a higher number of variables would be applicable in the binary models concerning the more frequent lexemes; however, it can be justified since it ensures that every individual binary model constituting the overall multinomial model will be generally valid and relevant. Furthermore, because of this limitation on the number of variables, there is no space for the consideration of interaction variables in the multivariate regression analysis.

The selection of variables that are actually included in the model, represented by $X = \{X_1, \dots, X_M\}$ in the formulas, is based on both the univariate results and the subsequent pairwise comparisons of the originally chosen contextual features, which should, on its part, be based first and foremost on domain-specific knowledge such as earlier studies and descriptions and theoretical motivations (Harrell 2001: 66). In general, the selected features should be both frequent enough and broadly distributed in order to rule out idiosyncratic associations, and therefore should at the least exceed

the minimum frequency threshold established in univariate analysis and should have observed occurrences with more than one of the studied lexemes. Furthermore, features which have in the pairwise analysis (using Cramér's V or PRE measures such as the Goodman-Kruskal τ or Theil's Uncertainty Coefficient U) been observed to correlate substantially either positively or negatively with each other, thus exhibiting *collinearity* or *multicollinearity*, should be considered carefully. The reason for this is that including both such features will not increase the explanatory power of the model, though it would, nonetheless, reduce the number of other features which can be included in the model, given the limiting sample size.

In addition, features which correlate with some other features, or groups of features perfectly, that is, features that could be categorically determined by some other individual feature (*exact collinearity*) or groups of features (*exact multicollinearity*), are troublesome as their inclusion in the model will not allow for the proper computation of the regression equation coefficients (Cohen et al. 2003: 419-430). In the case of a complementary distribution, exhibiting a perfect negative association, the solution is to include only one of the two features. The same also applies for cases of directional association arising from descriptive redundancy and overlap. Nevertheless, the reduction of highly correlating variables is not entirely unproblematic when there is no clear theoretical motivation for the selection of the variable(s) to be dropped, since the removal of a truly relevant variable will distort the estimates concerning the remaining variables (Cohen et al. 2003: 426-427).

A sophisticated method for substantially reducing collinearity is to use *Principal Components Analysis* (PCA) to transform the original variables into new aggregate variables based on the resultant principal components, and thereafter to undertake regression analysis with the transformed aggregate variables, discarding the smallest component(s) having the least variance in relation to the original variables, since the latter also account for most of the original collinearity. However, the resultant coefficients for the aggregate variables seldom have a clear meaning by themselves, and would have to be transformed back to the coefficients for the original variables. Furthermore, discarding the smallest components means that regression analysis with the aggregate variables is not equivalent to the results based on the original variables (Cohen et al. 2003: 428-429, see also Harrell 2001: 66-74, 75 [Figure 4.3] for other methods of variable reduction by clustering). However, in order to avoid adding to the complexity of this already multi-staged study, I will keep to working with only the original variables.

Nevertheless, when the number of nominal variables is quite large, it is probable that some intercorrelation remains among the features which can never be fully purged. Specifically, logically related, mutually exclusive feature sets such as the person/number, mood, or the semantic classifications of each syntactic argument type, corresponding to the binary dummy variables discussed below, are always partially correlated (Cohen et al. 2003: 311). It has also been observed that variables which correlate do not necessarily diminish the explanatory power of the model as much as one might expect, as long as the correlation is not limited only to the observed data but is sufficiently general to exist also in unseen, new data (Harrell 2001: 65). However, the role of a feature's significance or non-significance as observed in univariate analysis is of lesser importance. To the contrary, it has, in fact, been observed that leaving out features deemed insignificant in univariate analysis can

inflate and distort the weights and relationships of the remaining features (Harrell 2001: 56, 61; see also Bresnan et al. 2007). This is not to say that no superfluous features should be pruned, but neither should this practice be carried out to the extreme.

The individual binary feature variables can also be understood as the result of *dummy-variable coding* of variables with multiple classes (see Agresti 2002: 177-179; Cohen et al. 2003: 302-320; Harrell 2001: 14). In such a scheme, one reformulates each multi-class variable with c classes as $c-1$ dummy variables, with one of the classes, typically the most frequent or prototypical one chosen as the reference value for which all the $c-1$ dummy variables are FALSE (or =0). The reference class should not be an infrequent one, nor should it be a “waste-basket” or “dump” category. Though there are other binary coding alternatives for multi-class variables, the notion of *prototypicality* inherent in dummy-coding is quite appealing from the viewpoint of current linguistic theory. For instance, instead of a single multi-class variable for the person/number feature of the studied lexemes, which has $c=6$ values/classes corresponding to each of the theoretically possible six person number features, plus their nonoccurrence as in practice a seventh value, we can minimally have $c-1=5$ binary variables, each corresponding to one class of the person/number feature, with the THIRD PERSON SINGULAR as the reference class. This choice of the reference class can be based on the previous research indicating that the THIRD PERSON SINGULAR is, not only the most common person/number feature for any Finnish verb (Karlsson 1986; cf. Arppe 2006c), but this feature can, with justification, be considered the most prototypical if not the most natural one, too (Karlsson 1986: 26-27, in criticism of Mayerthaler 1981). The statistical motivation for only $c-1$ dummy variables and for not having a redundant dummy variable of its own for the reference class is that the redundancy, and other types of exact correlation with individual variables or variable groups, will not allow for the fitting of the regression equations uniquely.

However, this selection (and reduction) of reference classes for multi-class variables is problematic for such a large number of feature types considered in this study, because many of the multi-class variables are applicable for only a subset of the theoretically possible cases, and are thus not universally mutually exclusive. Furthermore, the linguistic descriptive system is not fully unambiguous. For instance, as the person-number features concern, strictly speaking, only ACTIVE FINITE forms, and by extension those NON-FINITE PARTICIPIAL forms which are used as CLAUSE-EQUIVALENTS and which can semantically have a person/number feature in the form of a possessive suffix (e.g., *harkittuaan* PCP2+ PASS+PTV+POSS:3 ‘once he has/had considered’), the reference class cannot be uniquely determined by the joint FALSE values of the dummy binary person/number features. This is due to the fact that these dummy-coded binary features are jointly FALSE even when none of the person/number features can be present, such as in PASSIVE forms. In addition, though the person/number features firstly concern only ACTIVE FINITE forms in contrast to PASSIVE forms, as was exemplified above, the PASSIVE feature may be associated with features semantically representing person/number in PARTICIPLE forms when used as CLAUSE-EQUIVALENTS. Furthermore, though we may designate quite easily a unique prototypical reference class for multi-class variables such as person/number or the semantic/structural type of AGENT, being the THIRD PERSON SINGULAR and the (human) INDIVIDUAL, respectively, this becomes more difficult for other syntactic arguments such as the PATIENT.

The more I consider this issue, the more I am inclined to believe that we cannot determine reference classes that would apply universally for all possible syntactic argument combinations with the studied lexemes. Instead, these reference classes are interrelated with each other, and some of them, either individually or in combinations, are particular to individual lexemes, for example, the *että*-clause as a PATIENT with *ajatella*. For these reasons, the focus in variable selection is on the identification of high mutual correlation among the binary variables, in addition to the identification of narrowly distributed variables. It is here that the results of both the univariate analyses and the bivariate comparisons, in combination with an overall linguistic perspective, are necessary. Potential reference classes of multi-class variables, such as the THIRD PERSON SINGULAR for person/number features and (HUMAN) INDIVIDUAL AGENTS are retained as variables, unless they are observed to be excessively intercorrelated with other variables. Nevertheless, there are (typically complementary) features, such as the FINITE/NON-FINITE distinction, which apply for the entire data. In their case, the number of binary variables included in the regression analysis must and will be reduced.

As a final note, it might be prohibitively difficult, for a variety of reasons, to adhere in later research to the full model and all of its explanatory variables which this study will build upon. For instance, only a subset of the explanatory variables might be readily coded in new corpus data, but, nevertheless, one might be interested in comparing the results of a smaller model with those presented here. Furthermore, the full model scrutinized in this study most probably is not the most parsimonious one, no matter how thoroughly it covers the studied phenomenon in detail. In such a situation one can consider the full model as a “gold standard”, against which one can then compare simpler models (Harrell 2001: 98-99). In order to facilitate such comparisons, I will therefore in this study also fit and test several simpler models with the same data as is used with the full model. These will include models containing 1) only node-specific morphological features, 2) verb-chain general morphological features, 3) syntactic argument types, *without* their semantic and structural classifications, 4) verb-chain general morphological features together with syntactic argument types without their subtypes, and 5) the aforementioned features and the most common semantic classifications of AGENTS and PATIENTS, with the less frequent types collapsed together, whenever possible.

3.4.3 Alternative heuristics of multinomial regression analysis

As was noted above, multinomial regression proper is based on selecting a baseline category among the outcomes. In the case of the four selected lexemes, this would undoubtedly be *ajatella*, as it is the most frequent and as it has the widest range of possible connotations presented earlier in Section 2.3.2. In the interpretation of the explanatory variables such a baseline setting is practical in contrasting the three other lexemes against this prototypical one. However, if we also would rather contrast *ajatella*, or in fact any individual of the four lexemes against the rest, and see which explanatory variables are distinctive, multinomial regression proper, assuming a baseline category, does not seem the most appropriate set-up. However, a number of heuristics, in addition to the multinomial base-line category model, have been

developed for analyzing such polytomous responses with logistic regression.⁶² These heuristics are all based on the splitting of the polytomous case into a set of dichotomous cases, for which the binary logistic regression model can then be applied separately, hence, they can be called *binarization techniques* (Fürnkranz 2002: 722-723). The differences among the heuristics are in the strategies according to which the decomposition into binary models and their overall fitting is undertaken. The relevant heuristics, in addition to the baseline category multinomial model already presented above, are 1) one-vs-rest classification, 2) pairwise classification, 3) nested dichotomies, and 4) ensembles of nested dichotomies. A concise presentation of all these and a few more heuristics can be found in Frank and Kramer (2004). In general, it has been observed that the process of separately fitting the binarized models does not generally have a substantial (detrimental or differentiating) effect on the overall results, in comparison to simultaneously fitting a proper multinomial model. Nevertheless, the latter is sometimes considered preferable as the most “elegant” solution (Hosmer and Lemeshow 2000: 277-278; Agresti 2002: 273-274).⁶³

One-vs-rest classification

The heuristic of *one-vs-rest* classification (e.g., Rifkin and Klautau 2004, also referred to as *one-vs-all*, *one-against-all*, *OVA*, or *unordered* classification) is based on contrasting (“singling-out”) each individual class k of the total of K outcomes against all the rest, with these $K-1$ classes lumped together into one alternative outcome. Thus, the one-vs-rest heuristic consists of K binary regression models, which are each trained with the entire data (see formulas 3.23-3.25 below). It is certainly conceptually simple, and according to Rifkin and Klautau (2004: 102), it has been independently discovered time and again by numerous researchers.⁶⁴

For the four studied lexemes, the exhaustive listing of the contrasts are therefore *ajatella* vs. *miettiinä* or *pohtia* or *harkita*, *miettiinä* vs. *ajatella|pohtia|harkita*, *pohtia* vs. *ajatella|miettiinä|harkita*, and *harkita* vs. *ajatella|miettiinä|pohtia*. In this setting, the regression coefficients of the individual binary models can be understood to highlight those feature variables which distinguish the individual outcome classes (i.e., lexemes) from all the rest, and they can meaningfully be studied together. A positive individual log-odds (coefficient) for some feature variable and the singled-out lexeme can be interpreted as the increased chances of the occurrence of this lexeme, when this particular feature is present in the context. In contrast, a negative log-odds would denote the decreased chances of the occurrence of this lexeme, translating into corresponding increased odds of any one of the three other lexemes occurring in such

⁶² Hereinafter, I will use *multinomial model* to refer to the heuristic where a set of (binary) baseline models are fitted simultaneously and in relation to each other with a given algorithm, often with the clarifying attribute “simultaneously-fitted” or “baseline-category” or “proper”. Using *polytomous models*, I will refer to the more general case of any heuristic for tackling polytomous outcomes which is based on logistic regression analysis, whether the component binary models are separately or simultaneously fitted.

⁶³ In fact, Fox (1997: 468, Note 34) does mention briefly that a symmetric alternative, with a probability estimation formula for each class k , is possible for the multinomial model, without the need for designating a baseline category. However, he notes that this would complicate the computations somewhat, and does not pursue it further.

⁶⁴ As a case in point, this was the heuristic we worked out together with Martti Vainio on our own, before scouring the literature and the Internet for alternatives.

a context. Consequently, if, in principle, a given feature has equal association with the singled-out lexeme and one but not all of the rest, since the other lexemes are lumped together, such a feature will not be treated as being as distinctive as it actually is.

Furthermore, given a particular constellation of values for the explanatory variables, the individual models yield direct probability estimates of the occurrence of the associated class k , or alternatively its nonoccurrence, implying the occurrence of any one of the $K-1$ complementary classes. In the prediction of the outcome, given a feature context X , the class for which the associated binary model yields the highest probability estimate wins, i.e., $\arg_k\{max_k[P_k(X)]\}$. As the binary logistic models are trained separately from each other, their joint probabilities are not necessarily exactly $\sum_{k=1...K}P_k(X)=1$. In fact, as a sneak preview of the multivariate results, for the total of 3404 instances in the data, the 95% CI of the instance-wise sums of probability estimates is $0.771 < \sum_{k=1...K}P_k(X) < 1.195$. In a sense, this could be interpreted as conformant with the *50/60 principle* concerning the acceptability ratings of alternative linguistic structures (see Arppe and Järvikivi 2007b).

(3.23) $P_k(X) = P(Y=k|X)$, with and $k=\{1, \dots, K\}$, and $P_{-k}(X) = P(Y=-k|X) = 1-P_k(X) = 1-P(Y=k|X)$ as the opposite case, i.e., the 'rest', so naturally $P_k(X) + P_{-k}(X) = 1$ for each binary model.

(3.24) $\log_e[P_k(X)] = \alpha_k + \beta_k X \Leftrightarrow P_k(X) = \exp(\alpha_k + \beta_k X)$

(3.25) $\beta_k X = \beta_{k,1}X_1 + \beta_{k,2}X_2 + \dots + \beta_{k,M}X_M$

with classes $k=\{1, \dots, K\}$, and M explanatory variables $X=\{X_1, \dots, X_M\}$, parameters $\beta = \{(\beta_{1,1}, \dots, \beta_{1,M}), (\beta_{2,1}, \dots, \beta_{2,M}), \dots, (\beta_{K,1}, \dots, \beta_{K,M})\}$, and constants $\alpha=\{\alpha_1, \dots, \alpha_K\}$

Pairwise classification

The heuristic of *pairwise* classification (e.g., Fürnkranz 2002, also referred to as the *round-robin*, *all-against-all*, *all-pairs*, and *AVA* classification) is based on the pairwise comparison of each class k_1 (of the altogether K classes) individually with every k_2 of the remaining $K-1$ classes with binary logistic models. In principle, the comparison of class k_1 against k_2 , i.e., $P_{k_1/k_2}(X)$ should be the mirror image of the comparison of class k_2 against k_1 , i.e., $P_{k_2/k_1}(X) = 1 - P_{k_1/k_2}(X)$. As a guarantee against this not always being the case in practice, for example, in computational implementations, the comparisons can be undertaken both ways, hence denoted as the *double-round-robin* technique. Thus, the pairwise heuristic amounts to as many as $K \cdot (K-1)$ binary logistic regression models, which are, however, trained with only the subset of the data having as the outcome one of the contrasted pair, $Y=\{k_1, k_2\}$, but none of the rest (see formulas 3.26-3.28 below).

For the studied four lexemes, there are in all $4 \cdot (3-1)=12$ contrasts, starting with *ajatella* vs. *miettiä*, *ajatella* vs. *pohtia*, *ajatella* vs. *harkita*, followed by *miettiä* vs. *ajatella*, and so forth. In this setting, the regression coefficients can be understood to highlight those features which distinguish the individual contrasted pairs from each other. Therefore, they do not have a direct overall interpretation such as the coefficients of the individual models in the one-vs-rest heuristic, even more so as the

binary models are trained with only the two contrasted lexemes at a time. Nevertheless, the pairwise odds derivable from the coefficients of the $K-1$ contrasts of each lexeme against the rest can be pooled for each lexeme by averaging them geometrically to provide a conservative approximate overall odds of each feature per lexeme (this geometric average of the odds-ratios corresponds to the arithmetic average of the log-odds, that is, coefficients, see formula 3.29). However, this method of aggregation may not perform satisfactorily in the contradictory case of one lexeme contrasting positively with another lexeme and negatively with a third lexeme.

In the prediction of outcome for a given context and constellation of features, direct probability estimates for each lexeme are not available, either. Instead, a voting scheme is used to aggregate the binary comparisons, where in its simplest (unweighted) form a lexeme k_1 gets one vote for each of its contrasted binary models for which its probability is $P_{k_1/k_2}(X) > 0.5$, given the context; otherwise, the vote goes to the contrasted lexeme k_2 instead, that is, when $P_{k_1/k_2}(X) \leq 0.5$. The lexeme k receiving the highest number of votes wins; in the case of a tie, the more frequent lexeme is selected.⁶⁵ Nevertheless, the number of votes per lexeme can be divided by the overall number of votes to produce a *very* rough approximation of the lexeme-wise probabilities, given a particular context (3.30). In principle, this setting with binary comparisons should produce better results in prediction when crucial distinctions are to be found between two individual lexemes instead of between one individual lexeme and all the rest. Furthermore, pairwise contrasting should be theoretically simpler in terms of the pairwise decision boundaries (Fürnkranz 2002: 724), but it remains to be seen what the actual effects are in the linguistic setting at hand.

$$(3.26) P_{k_1/k_2}(X) = [P(Y=k_1|X) | Y=\{k_1, k_2\}], \text{ and } P_{k_2/k_1}(X) = 1 - P_{k_1/k_2}(X) = 1 - [P(Y=k_1|X) | Y=\{k_1, k_2\}]$$

$$(3.27) \log_e[P_{k_1/k_2}(X) | Y=\{k_1, k_2\}] = \alpha_{k_1/k_2} + \beta_{k_1/k_2} X$$

$$(3.28) \beta_{k_1/k_2} X = \beta_{k_1/k_2,1} X_1 + \beta_{k_1/k_2,2} X_2 + \dots + \beta_{k_1/k_2,M} X_M$$

$$(3.29) \beta_{k_1,m} \approx (\beta_{k_1/k_2,m} + \beta_{k_1/k_3,m} + \dots + \beta_{k_1/K,m}) / (K-1), \text{ since the geometric average of the binary log-odds is } [e^{\beta^{(1)}} \cdot e^{\beta^{(2)}} \cdot \dots \cdot e^{\beta^{(K-1)}}]^{1/(K-1)} = e^{[\beta^{(1)} + \beta^{(2)} + \dots + \beta^{(K-1)}] / (K-1)}$$

$$(3.30) P_{k_1}(X) \approx \{n[P_{k_1/k_2}(X) > 0.5] + n[P_{k_2/k_1}(X) \leq 0.5]\} / [K \cdot (K-1)]; \text{ N.B. } 0 \leq P_{k_1}(X) \leq 0.5$$

with classes $k_1 = \{1, \dots, K\}$, and $k_2 = \{1, \dots, K\}$, with $k_1 \neq k_2$, and M explanatory variables $X = \{X_1, \dots, X_M\}$, parameters $\beta = \{(\beta_{1/2,1}, \dots, \beta_{1,M}), \dots, (\beta_{1/K,1}, \dots, \beta_{1/K,M}), (\beta_{2/1,1}, \dots, \beta_{2/1,M}), \dots, (\beta_{2/K,1}, \dots, \beta_{2/K,M}), \dots, (\beta_{K/1,1}, \dots, \beta_{K/1,M}), \dots, (\beta_{K/K-1,1}, \dots, \beta_{K/K-1,M})\}$, and constants $\alpha = \{\alpha_{k_1/k_2}, \alpha_{k_1/k_3}, \dots, \alpha_{K/K-2}, \alpha_{K/K-1}\}$

⁶⁵ As Fürnkranz (2002: 725, 738-739) concedes, this simplest possible voting procedure is most certainly suboptimal, but I will adhere to it in this study in order to avoid excessive additional complexity.

Nested dichotomies

In the technique of *nested dichotomies* (Fox 1997: 472-475; see also Cohen et al. 2003: 520-522; Frank and Kramer 2004), the original multi-class setting with K classes is recursively split into two subsets until there are only unary or binary subsets left, the whole of which can be represented as a binary decision tree of dichotomous contrasts. For any number of classes greater than two, there is always more than one way to split the classes,⁶⁶ and the total number of these possible partitions grows extremely quickly with the number of classes, according to the recursive formula $T(K)=(2 \cdot K-3) \cdot T(K-1)$, where $T(1)=1$.⁶⁷ In contrast, the number of binary models for an individual partition is quite moderate at $K-1$ (which would each be trained with the subset of the data relevant to each partition as in the pairwise heuristic).

The four studied lexemes could be partitioned in $T(4)=15$ ways, such as $\{\textit{ajatella}$ vs. $\{\textit{miettiin} \text{ vs. } \{\textit{pohtia}$ vs. $\textit{harkita}\}\}$, or $\{\{\textit{ajatella}$ vs. $\textit{miettiin}\}$ vs. $\{\textit{pohtia}$ vs. $\textit{harkita}\}\}$, each involving $4-1=3$ binary models. However, nested dichotomies are recommended only when some particular partition can be motivated over the rest on the basis of domain-specific knowledge (Fox 1997: 472). As the studied lexemes already belong to a semantically tightly-knit synonym group, at least to my mind there is no obvious single partition that could be argued to be above the rest on linguistic grounds. For instance, one could envisage contrasting the most frequent and semantically broadest *ajatella* against the rest, or one could consider grouping the etymologically agriculture-originated *pohtia* and *harkita* against the more neutral *ajatella* and *miettiin*. But, one could just as well differentiate *harkita* from the rest on the basis of Pajunen's (2001) classification presented above in Section 2.3.1.

Nevertheless, nested dichotomies have the attractive characteristic that this heuristic allows for the straight-forward calculation of probability estimates for the individual classes – without approximations and post-processing. These are calculated simply by multiplying the probabilities on the path from the root through the relevant internal nodes to each particular leaf (i.e., lexeme) of the binary classification tree. In the case of the partition $\{\textit{ajatella}$ vs. $\{\textit{miettiin}$ vs. $\{\textit{pohtia}$ vs. $\textit{harkita}\}\}$, the probability of the outcome $Y=\textit{harkita}$ for a given context and features (represented as X) would thus be $P(Y=\{\textit{miettiin}, \textit{pohtia}, \textit{harkita}\}|X) \cdot P(Y=\{\textit{pohtia}, \textit{harkita}\}|X) \cdot P(Y=\{\textit{harkita}\}|X)$ (for an exact formula see Kramer and Frank 2004). However, the existing literature does not explicitly present a method for aggregating lexeme-specific estimates of the related odds-ratios of the feature variables, which are of specific interest in this linguistic study. Nevertheless, the probability structure of the partitioning would suggest, as one possible avenue of aggregation, that one would multiply the relevant sequences of odds-ratios from the root to the lexeme, in a fashion similar to the computation of the probability estimates.

⁶⁶ In the simplest case, $T(n=3)=3$ with $\{A, B, C\} \rightarrow \{\{A, B\}, C\}$ or $\{A, \{B, C\}\}$ or $\{\{A, C\}, B\}$

⁶⁷ $T(1)=1$; $T(2)=1$; $T(3)=3$; $T(4)=15$; $T(5)=105$, and so forth.

Ensembles of nested dichotomies (ENDs)

As a solution to the theoretical problems of selecting one single nested partition over the rest, Frank and Kramer (2004) propose using an *ensemble of nested dichotomies* (denoted by the acronym *END*). Their line of argumentation in this is that, when none of the individual partitions can theoretically be established as substantially better than the rest, it would make sense to regard each partition tree as equally likely and, therefore, to study their overall behavior as an *ensemble*, hence the name. The probability estimates of individual outcome classes can then be calculated as averages of the estimates derived from the individual partitions. As the number of binary models necessary in all partitions grows even faster than the number of partitions, amounting to $[3^K - (2^{K+1} - 1)]/2$ of individual binary models for K outcome classes, their number has to be restricted in some manner. For this purpose, Frank and Kramer show that using a random selection of 20 partitions (with $K-1$ binary models for each partition) is sufficient in most cases for achieving “close-to-optimum” performance. In the case of the four lexemes studied here, however, this would not make any difference as the overall number of partitions is $T(4)=15 < 20$.

Nonetheless, if we had included only one more lexeme in the studied synonym group, the overall number of possible partitions would have risen to $T(5)=105$, and with six lexemes the figure would have continued to rise exponentially to $T(6)=945$. In such cases, a smaller, randomly sampled set of partitions might be more desirable in order to decrease the computational load, specifically in the resampling schemes to be discussed later. Furthermore, the approximation of lexeme-specific odds-ratios of the individual feature variables would be complicated even further, as one would then have to take all the different partitions into account. In principle, however, these could be calculated as the averages of the partition-specific aggregated odds-ratios, which for each partition would be, in turn, the products of the relevant sequences of odds-ratios from the root to the lexeme.

Comparing the heuristics and their characteristics

Table 3.40 presents a comparison of the characteristics of the various heuristics for polytomous logistic regression presented above, and thus also their pros and cons from the perspective of this linguistic study. In order to obtain both lexeme-specific parameters for the contextual features, without having to select one lexeme as a baseline, and probability estimates for the occurrences of each lexeme, the one-vs-rest heuristic is the most appealing of those available. To its benefit, it is also methodologically simple, as both the parameters and the probability estimates are directly derived from the binary logistic regression models of which it consists. In contrast to the pairwise heuristic that I tentatively applied in Arppe (2006c), the one-vs-rest heuristic requires considerably fewer binary logistic models (in this case 4 vs. 12, and this ratio becomes increasingly better with the growth of the number of outcomes). In addition, the one-vs-rest heuristic provides the parameters lexeme-wise directly without any additional and approximate aggregation from the pairwise contrasted models. Furthermore, Rifkin and Klautau (2004: 102) argue forcefully that, contrary to the commonly held assumption, one-vs-rest is not less accurate than other, typically more sophisticated heuristics.

Nevertheless, I will compare all the different heuristics with respect to their prediction accuracy, as that is the purpose which they seem most geared towards, and since it will also give some indication of whether the parameters of the underlying binary models might be worth further investigation. What is more, the underlying concept of nested dichotomies and ENDS is appealing from the linguistic viewpoint, since I consider it conceivable that each possible partition would represent different perspectives in the contextual behavior of the studied lexemes. For instance, one partition might concern the types of AGENTS that the studied lexemes prefer, another the type of PATIENTS they occur with, a third the types of person/number they appear in, and so on. Along this line of thinking, an ensemble of these partitions would then reflect the aggregated effect of these different types of contextual features in the selection of synonymous lexemes. Furthermore, comparing the different nested partitions could be used to study how the studied lexemes relate to each other; if the predictive capabilities of a given partition were observed to be significantly better than that of the others, one could consider the partition in question to best represent the structure of studied lexemes as group.

Table 3.40. The general characteristics and pros and cons of various methods/heuristics for polytomous regression.

Heuristic/ characteristics	Multinomial (baseline category)	One-vs- rest	Pairwise	Nested dichotomy	Ensemble of nested dichotomies
Number of constituent binary models	$n_{lex}-1$	n_{lex}	$n_{lex} \cdot (n_{lex}-1)/2$ (round-robin) $n_{lex} \cdot (n_{lex}-1)$ (double-round- robin)	$n_{lex}-1$	~ 20 partitions (each with $n_{lex}-1$)
Lexeme- specific odds- ratios for feature variables	No (Every lexeme against the baseline)	Yes (Every lexeme against the rest)	No (Approximation by geometric averages of binary odds- ratios)	Yes (Products of binary odds-ratios)	Yes (Averages of products of binary odds-ratios)
Probability estimates for lexemes (i.e., outcomes)	Direct	Direct $P_{lex/rest}(X)$	No	Direct (Product of probabilities at nodes in partition tree)	Direct (Average of products of probabilities at nodes in partition tree)
Selection of lexeme in prediction	Probability- based $\arg_{lex} \max(P_{lex X})$	Probability- based $\arg_{lex} \max(P_{lex X})$	Voting $\arg_{lex} \max$ $\{n[P_{lex1/lex2}(X)>0.5]$ + $n[P_{lex2/lex1}(X)\leq 0.5]\}$	Probability- based $\arg_{lex} \max(P_{lex X})$	Probability- based $\arg_{lex} \max(P_{lex X})$
Other	Necessity of baseline category	May not discover pairwise distinctions	May exaggerate pairwise distinctions, and the behavior with contradictory distinctions is problematic	Selection of a single appropriate partition may be difficult or impossible	-

3.4.4 Evaluating the polytomous logistic regression models and their performance

There are several perspectives along which polytomous logistic regression models with categorical explanatory variables can be evaluated. First of all, analogously with “ordinary” linear regression, we can assess to what extent overall the logistic models fit and account for the data on which they are based. Secondly, we can test how well the models generalize and how accurately they are able to predict outcomes with new, unseen data, with which they have not been fitted and trained. As a variant of this, we can also test how well the models can predict the outcomes in the data that they were originally trained on. Thirdly and finally, we can use various resampling schemes to evaluate both the accuracy of prediction and the robustness of the effects represented by the estimated parameter coefficients of the explanatory variables in the models. Since logistic regression models in the first place estimate the *probabilities* of occurrence and not the categorical *occurrences* or nonoccurrences of alternative choices, in principle the assessment of the fit of a model with the original data should take precedence over the evaluation of the model’s prediction accuracy (Harrell 2001: 248-249). Focusing primarily on prediction accuracy is justified when classification is an explicit goal; otherwise, it should only be considered a supplementary form of evaluation (Hosmer and Lemeshow 2000: 160).

Evaluation of model fit with original data

The evaluation of the overall fit of logistic regression models is based on the measure of their *decreased deviance* (Agresti 2002: 139-142 and 186-187; Cohen et al. 2003: 499-506; Fox 1997: 450-451), in contrast with the *increase* of explained *variance* as observed in conjunction with “ordinary” linear regression models. Deviance (denoted as D) is a relative measure, and it is based on the *lack-of-fit* of a given model M_1 compared to another, typically simpler or baseline, model M_0 . This lack of model fit is represented by their associated *maximum likelihoods* L_1 and L_0 , in terms of which deviance is defined as the natural logarithm of their ratio, denoted as LR (formula 3.32). In turn, likelihood L is the joint probability of the actually observed outcomes, as assigned by any particular model given the contextual data with which it is fitted. As this joint probability is the product of the individual probabilities, for reasons of simpler calculation, the logarithm of the likelihood, that is, *log-likelihood*, is studied instead (formula 3.31, see Eliason 1993: 7-8, equations 1.6 and 1.7). One should note that in a polytomous case, with multiple possible outcome classes, only the probability corresponding to each actually observed outcome (and its particular context) is included in the calculation of likelihood. The probabilities corresponding to the non-observed outcomes are not considered for any instance, though these other outcomes may, in principle, be possible and perhaps be associated with a substantial probability estimate.

(3.31) $L = \prod_{i=1 \dots N} P(Y_i) \Leftrightarrow \log_e L = \sum_{i=1 \dots N} \log_e [P(Y_i)]$, for $i = \{1, \dots, N\}$ originally observed outcomes, $Y_i = \{Y_1, Y_2, \dots, Y_N\}$, each outcome belonging to one class k of altogether K classes, i.e., $\forall Y_i \in \{1, \dots, K\}$.

(3.32) $D = -2\log_e(L_1/L_0) = -2[\log_e(L_1) - \log_e(L_0)] = -2\log_e(LR)$.

The maximum likelihood for any sample of data theoretically has two extreme end-points between which it can vary, these being perfect maximum likelihood $L_{perfect}$ and null maximum likelihood L_{null} . *Perfect maximum likelihood* for some data is by definition equal to 1.0, and would in principle be attainable with a perfect model (sometimes also called a *saturated* model, e.g., Agresti 2002: 187). In such a perfect model, each observed outcome would be matched by an explanatory variable of their own, resulting in $P(Y_i=k|X_i)=1$ always and only when $Y_i=k$, and $P(Y_i=k|X_i)=0$ otherwise when $Y_i \neq k$ (for each outcome k in $\{1...K\}$). In contrast, *null maximum likelihood* for the same data is the (almost) opposite case where the model would be null and would consist of only an intercept and no explanatory variables at all. In a dichotomous case, the intercept is exactly the log-odds of the outcome belonging to class k instead of the other, complementary class, $\alpha = \log_e[(n_k/N)/(1 - n_k/N)]$, i.e., $\text{logit}[n_k/N]$ (see Harrell 2001: 228), leading to the corresponding null log-likelihood in 3.33. With polytomous outcomes, as is the case here, the intercepts associated with the null maximum likelihood are the logarithms of the overall probabilities for each individual class k , i.e., $\alpha_k = \log_e(n_k/N)$ (see Menard 1995: 84), which we could expect without knowledge of the influence of any explanatory variables included in a model. This yields the null maximum log-likelihood presented in 3.34.

For the aforementioned extreme ends of likelihood we can calculate their mutual difference ratio, designated as *null deviance* D_{null} and presented in 3.35, which is also the maximum deviance that any model could theoretically account for the given sample of data. Then, for a model with m explanatory variables, we can also calculate its deviance D_{model} in relation to the perfect and the null cases (formulas 3.36 and 3.37). The maximum likelihood values for a sample of training data and a particular logistic model are estimated as a part of the iterative algorithm through which this model and its coefficients are fitted with the data, with the goal to maximize the associated overall likelihood for the model and the data sample (Cohen et al. 2003: 498-499). Thus, once we have a fitted model thanks to some statistical software, we can calculate the associated maximum log-likelihood simply by adding up the logarithms of the probabilities estimated by the fitted model (or their combinations) for each originally observed outcome. Knowing the explanatory variables, the expectation naturally is that, the estimated probabilities and thus also the likelihood (as well as the log-likelihood) would be greater *overall* than the simple overall probabilities of the classes alone, though for some *individual* cases the estimated probabilities might actually turn out to be less than expected at the null level.

$$(3.33) \log_e L_{null, dichotomous} = n_k \cdot \log_e(n_k) + (N - n_k) \cdot \log_e(N - n_k) - N \cdot \log_e(N)$$

$$(3.34) \log_e L_{null, multinomial} = \sum_{k=1...K} \{n_k \cdot \log_e[P(Y_i=k)]\} = \sum_{k=1...K} [n_k \cdot \log_e(n_k/N)]$$

$$(3.35) D_{null} = -2[\log_e(L_{null}) - \log_e(L_{perfect})] = -2[\log_e(L_{null}) - \log_e(1)] = -2\log_e(L_{null})$$

Where n_k is the total number of outcomes for class k so that $Y_i=k$ and N is the total sample size (see Harrell 2001: 228, equation 10.24; Menard 1995: 84)

$$(3.36) \log_e L_{model} = \sum_{i=1...N} \log_e[P(Y_i=k|X_i)]$$

Where N is the total sample size, Y_i is the original i th outcome, each with altogether K possibilities, so that each $\forall Y_i \in \{1, \dots, K\}$, and $P(Y_i=k|X_i)$ is the fitted probability

estimate for context X_i corresponding to the actually observed outcome $Y_i=k$. Therefore, the estimated probabilities for any other possible outcome, i.e., $P(Y_i \neq k | X_i)$, are not considered in the calculation of the overall likelihood (cf. Eliason 1993: 7-8).

$$(3.37) D_{\text{model}} = -2[\log_e(L_{\text{model}}) - \log_e(L_{\text{perfect}})] = -2[\log_e(L_{\text{model}}) - \log_e(1)] = -2\log_e(L_{\text{model}})$$

The purpose of all the above formulations is to lay the ground for a measure of evaluating how much of the overall deviance a particular model that we have selected can account for. For this purpose, there are a variety of formulas available, but as none of them have been shown to be clearly superior to the rest, I will settle on the simplest one, that is, R_L^2 (formula 3.38), presented by Hosmer and Lemeshow (1989: 148; see also Menard 1995: 19-24; Fox 1997: 450-451; Hosmer and Lemeshow 2000: 165-166; Cohen et al. 2003: 502-504). The R_L^2 measure is analogous in structure with the multiple-correlation coefficient R^2 used in ordinary linear regression, but based on deviance as defined above it should not be confused with the proportion of variance in the data that the model is able to account for.⁶⁸ Furthermore, one should note that, despite this structural similarity, all of the logistic R_L^2 measures yield values which are typically quite low in comparison to the those encountered in the evaluation of linear regression models, even when they might represent the data accurately (Hosmer and Lemeshow 2000: 167). Finally, it is not uncommon to see the D_{model} measure used as the basis for testing the goodness-of-fit of the associated logistic regression model, by considering this deviance as asymptotically χ^2 -distributed, with $df=N-(m-1)$. However, since there is controversy as to whether this practice is in fact at all justified (e.g., Cohen et al. 2003: 504-506; for criticism, see Harrell 2001: 231 and also Baayen 2008: 217-218), I will not pursue that line of evaluation further here.

$$(3.38) R_L^2 = (D_{\text{null}} - D_{\text{model}}) / D_{\text{null}} = 1 - D_{\text{model}} / D_{\text{null}} = (\log_e L_{\text{model}} - \log_e L_{\text{null}}) / \log_e L_{\text{model}} \\ = 1 - \log_e L_{\text{null}} / \log_e L_{\text{model}}$$

One should note that when we use a heuristic based on a set of separately trained binary logistic regression models in order to accomplish polytomous logistic regression, the individual binary models are not fitted by maximizing the (log-)likelihood of all the polytomous outcomes, but only those binary outcomes at a time which are considered in the individual models. In fact, in such a case Hosmer and Lemeshow (2000: 280-281) suggest studying first the individual fits of the set of binary models, and then making a descriptive assessment of the overall fit on the basis of the component results. Nevertheless, as long as the heuristic produces direct probability estimates for all outcomes and classes, albeit via component models or their combinations, we can calculate an overall estimate of their (log-)likelihood and deviance and thus evaluate the overall fit of the multiple binary models considered together. In doing this, as specifically in the case of the *one-vs-rest* heuristic the sum of the probabilities of the binary models is not necessarily exactly equal to 1.0, the probability estimates should probably be scaled to take this possible variation into

⁶⁸ These log-likelihood based measures such as R_L^2 are sometimes characterized as *not* rendering themselves to *intuitively* easy and natural interpretation, as they do not correspond to the R^2 measures of linear regression in representing directly explained variance, and are thus, in the views of some, not to be recommended (e.g., Hosmer and Lemeshow 2000: 164, referring to criteria originally presented by Kvalseth in 1985). To the contrary, I find their basis in the overall outcome probabilities as an attractive one, as that is exactly what logistic regression purports to model, and thus, to my mind, they are not at all that obscure.

account. However, since the *pairwise* heuristic hardly provides any direct estimates of overall probability for the polytomous outcomes and because transforming the votes into probabilities are, at best, only coarse approximations, in its case it does not make much sense to evaluate the overall fit of the combination of the constituent binary models in terms of deviance and log-likelihood as presented here above. Furthermore, in the case of the *nested dichotomies*, the overall deviance can alternatively be calculated simply as the sums of deviances (based on the maximum log-likelihoods) of the individual binary models determined by the partition, due to their mutual independence (Fox 1997: 473-474).

Evaluation of prediction efficiency and accuracy

As much as the evaluation of the log-likelihood and deviance of the selected model, implemented according to the different heuristics, would in principle be the most appropriate way to evaluate the fit of the (polytomous) logistic regression model and its constituent set of explanatory variables with the data, it may be more valuable for the comparison of different selections of explanatory variables than for the overall evaluation of the model (see Hosmer and Lemeshow 2000: 167; Agresti 2002: 186-187). This is even more so as the R_L^2 –as well as the other “pseudo”- R^2 measures – do not have a natural interpretation as such; they rather indicate whether or not a model with its associated explanatory variables is better than another. In contrast, the ability of the model and the heuristic it has been implemented with to make predictions about which lexeme will occur in a given context is immediately more understandable, thus increasing the worth of prediction accuracy as an evaluation method of polytomous regression models. What is more, many of the heuristics presented above, namely, the one-vs-rest and pairwise classification, not to mention ensembles of nested dichotomies, have been developed with classification clearly in mind, which is evident in how these heuristics are presented and evaluated against other alternatives. It is in such a case that Hosmer and Lemeshow (2000: 167) regard the evaluation of classification accuracy as also appropriate.

We should remember, however, that classification as a task is categorical in nature and that it masks the underlying probabilities, especially in a polytomous setting with more than two alternatives; of the four studied lexemes, one class k can be selected as well as any other with a probability of just over $P(Y_i=k|X_i) > 1.0/4 > 0.25$, if the other three are only slightly less (and approximately equally) probable, as with an overwhelming preference represented by, e.g., $P(Y_i=k|X_i) = 0.9$ (cf. the thorough discussion regarding binary outcomes and probabilities in Hosmer and Lemeshow 2000: 156-160). Furthermore, from the linguistic perspective, since we are dealing with a synonymous set of lexemes, we may expect relatively similar underlying probabilities instead of significant dispersion in their values, as, in principle, on the basis of the previous descriptions presented above in Section 2.3.2, any individual one of the four lexemes can be used in most, if not all of the studied contexts. More specifically, if some context allows for genuine linguistic variation, at least as defined by the selected feature variables in the model, categorically selecting always one lexeme over the others on the basis of possibly a very small difference in estimated probabilities would not properly reflect the reality of linguistic usage.

By way of illustration, if for some fixed set of explanatory values X_i for a recurrent context a lexeme k receives the probability estimate $P(Y_i=k|X_i)=0.51$, which is thus the maximum value for this context, and there are exactly 100 instances of such a specific context in the original data, this means that we could expect lexeme k to have occurred $0.51 \cdot 100 = 51$ times, and any one of the three other lexemes the remaining $100 - 51 = 49$ times, each with their individual probability estimates corresponding to their proportions in the original data. However, a prototypical classification rule $Y_i=k \Leftarrow P(Y_i=k|X_i) > 0.50$ (or $Y_i=k \Leftarrow \arg_k \max [P(Y_i=k|X_i)]$) would result in lexeme k being predicted to occur for every instance of the specific context, in this case 100 times out of the 100. This clearly does not reflect the distributions and associated proportions of occurrence in the original data. In this respect, the scrutiny of the entire probability distributions for all the polytomous outcome classes retains an important role.

Our expectations concerning the prediction of outcome classes can, in fact, be divided into two types, namely, classification and prediction models,⁶⁹ which have an effect on how the efficiency and accuracy of prediction is exactly measured (Menard 1995: 24-26). In a pure *prediction model*, we set no *a priori* expectation or constraint on the overall frequencies of the predicted classes. Indeed, it would be acceptable for all predictions to be fully homogeneous and belong to only one single class, even though the training data may have contained (many, or at least some) occurrences of other classes. To the contrary, in a *classification model* our expectation is that the predicted outcome classes will, in the long run, end up having the same proportions as are evident in the training data. That is, we *a priori* expect heterogeneity among the predicted outcomes. The complete homogeneity of predicted outcomes would entail failure for a classification model, whereas it would be an acceptable result for a prediction model.

In this linguistic study, the prediction model entails our acceptance of the possibility that the selected lexeme – in any context in question – would be one and the same, this lexeme probably being the most frequent one, that is *ajatella*. In principle, we would then regard the four lexemes as absolute synonyms, fully interchangeable with each other in all possible contexts. However, the classification model entails that we expect, firstly, all four of the studied THINK lexemes to turn up as predicted outcomes and, secondly, with similar proportions as were observed in the original data. In this case, we assume that the lexemes do have minute semantic differences, which should become evident through their (at least slightly) different contexts of usage, that is, the four lexemes are only near synonyms. The classification model is more difficult to satisfy, especially because classification schemes tend to favor the most frequent class (Hosmer and Lemeshow 2000: 157). Nonetheless, it is also more in line with the views of current lexicographical theory and with what the original data suggests about the four synonyms.

⁶⁹ In fact, Menard (1995) presents also a third type of prediction model, namely, for *selection* (with some *a priori* fixed *selection ratio*), but only the two types discussed here are directly applicable – without any potential need for possible adjustments – to the classification tables in this study.

Table 3.41. Prediction and classification table n for the studied four THINK lexemes; $n_{1,1}$ corresponds to $\sum(\text{Predicted}=\text{ajatella} \wedge \text{Original}=\text{ajatella})$, $n_{1,2}$ corresponds to $\sum(\text{Predicted}=\text{miettiä} \wedge \text{Original}=\text{ajatella})$, $n_{2,1}$ corresponds to $\sum(\text{Predicted}=\text{ajatella} \wedge \text{Original}=\text{miettiä})$, and so forth.

Original/Predicted	ajatella	miettiä	pohtia	harkita	$\sum(\text{Original})$
ajatella	$n_{1,1}$	$n_{1,2}$	$n_{1,3}$	$n_{1,4}$	$\sum n_{1\cdot}$
miettiä	$n_{2,1}$	$n_{2,2}$	$n_{2,3}$	$n_{2,4}$	$\sum n_{2\cdot}$
pohtia	$n_{3,1}$	$n_{3,2}$	$n_{3,3}$	$n_{3,4}$	$\sum n_{3\cdot}$
harkita	$n_{4,1}$	$n_{4,2}$	$n_{4,3}$	$n_{4,4}$	$\sum n_{4\cdot}$
$\sum(\text{Predicted})$	$\sum n_{\cdot 1}$	$\sum n_{\cdot 2}$	$\sum n_{\cdot 3}$	$\sum n_{\cdot 4}$	N

The starting point for evaluating prediction efficiency is to compile a prediction/classification table n , which is naturally a square matrix with the dimensions $K \times K$ in accordance with the number of original classes K . For each original class k , one then proceeds to count the distribution of the predicted classes (Table 3.41), with the original classes here being the row variable and the predicted classes the column variable. Frequency counts on the diagonal in the table indicate correctly predicted and classified cases, whereas all counts off the diagonal are incorrect. In addition to the correct predictions, one can in a polytomous setting, as is the case here, also directly scrutinize the prediction table with respect to how the incorrect predictions are distributed for each original outcome class. This is motivated by the fact that the degree to which two classes are getting mixed up can be seen as representative of the extent of their similarity in terms of the explanatory variables, that is, the similarity of lexemes as to their feature contexts in this study. Furthermore, for each class individually and for the classes overall, we can divide the predicted classifications into the four types presented in Table 3.42 and in formulas 3.39–3.42, on which the basic measures of prediction efficiency are based.

Table 3.42. The four different classes of predictions

Original/Predicted	Class	\neg Class (=Other)
Class	TP ~ True Positive (=correct)	FN ~ False Negative (=incorrect)
\neg Class (=Other)	FP ~ False Positive (=incorrect)	TN ~ True Negative (=correct)

$$(3.39) \text{TP}(\text{class}=k) = n_{k,k}$$

$$(3.40) \text{FP}(\text{class}=k) = \sum_{i=1 \dots K} n_{i,k} - n_{k,k}$$

$$(3.41) \text{TN}(\text{class}=k) = N - \sum_{i=1 \dots K} n_{k,i} - \sum n_{i,k} + n_{k,k}$$

$$(3.42) \text{FN}(\text{class}=k) = \sum_{i=1 \dots K} n_{k,i} - n_{k,k}$$

Since this study concerns polytomous outcome cases where the models by design always have to select an outcome from the original cases, there exists no “extra” or non-classified category overall which should/could be classified as such and thus rejected. Such rejected non-cases will always belong to one of the other possible classes, and will thus not “fall out” of the classification scheme. In this respect, in the subsequent evaluation of prediction and classification efficiency the concept pair of *Recall* and *Precision* (Manning and Schütze 1999: 267-271), familiar from computational linguistics, feels most appropriate, as their computation in this

classification scheme makes more sense both class-wise and overall than the often used distinctive pairings of *Sensitivity* and *Specificity* (e.g., Cohen et al. 2003: 316). *Recall* is the proportion of original occurrences of some particular class for which the prediction is correct (formula 3.43, see Manning and Schütze 1999: 269, formula 8.4), whereas *Precision* is the proportion of all the predictions of some particular class, which turn out to be correct (formula 3.44, see Manning and Schütze 1999: 268, formula 8.3).

Sensitivity is, in fact, exactly equal to *Recall*, whereas *Specificity*, understood as the proportion of non-cases correctly predicted or classified as non-cases, that is, rejected (formula 3.45), is not really applicable in this study. The reason for this is that this correct rejection would translate into the (correct or incorrect) classification of the other lexemes as such, for the aforementioned reasons of mutually exclusive selection, and these non-cases would thus in fact partially overlap for the lexeme set as a whole, making its calculation for the classes overall pointless. Due to the same reasons in this classification scheme, *Recall* is equal to *Precision* for all the classes considered together (formula 3.46). Furthermore, there is a third pair of evaluation measures that one could also calculate, namely, *Accuracy* and *Error* (formulas 3.47 and 3.48); however, these are in general less sensitive than *Recall* and *Precision* to the class-specific counts (*True Positives*, *False Positives*, and *False Negatives*) which we are usually most interested in (Manning and Schütze 1999: 269-270). In a polytomous setting, their calculation class-wise would make little sense as the correct classifications of the class of interest are lumped together with the correct rejections of the other classes, while no attention is paid to whether the rejections of these other classes are indeed correctly classified. Taking all the above into account, in the actual evaluations of the prediction efficiency of the polytomous regression models only *Recall* and *Precision* will be calculated for each lexeme outcome, as well as overall *Recall*.

$$(3.43) \text{ Recall}_{\text{class}=k} = \text{TP} / (\text{TP} + \text{FN}) = n_{k,k} / \sum_{i=1 \dots K} n_{k,i} (= \text{Sensitivity}_{\text{class}=k})$$

$$(3.44) \text{ Precision}_{\text{class}=k} = \text{TP} / (\text{TP} + \text{FP}) = n_{k,k} / \sum_{i=1 \dots K} n_{i,k}$$

$$(3.45) \text{ Specificity}_{\text{class}=k} = \text{TN} / (\text{TN} + \text{FN}) \\ = (N - \sum_{i=1 \dots K} n_{k,i} - \sum_{i=1 \dots K} n_{i,k} + n_{k,k}) / (N - \sum_{i=1 \dots K} n_{i,k})$$

$$(3.46) \text{ Recall}_{\text{class}=1 \dots K} = \sum_{k=1 \dots K} n_{k,k} / N = \text{diag}(n) / N = \text{Precision}_{\text{class}=1 \dots K}$$

$$(3.47) \text{ Accuracy}_{\text{class}=1 \dots K} = (\text{TP} + \text{TN}) / N = \text{diag}(n) / N = \sum_{k=1 \dots K} n_{k,k} / N$$

$$(3.48) \text{ Error}_{\text{class}=1 \dots K} = (\text{FP} + \text{FN}) / N = [N - \text{diag}(n)] / N = 1 - \text{Accuracy}_{\text{class}=1 \dots K}$$

However, these aforementioned general measures do not in any way take into consideration whether prediction and classification according to a model, with the help of explanatory variables, performs any better than knowing the overall proportions of the outcome classes, corresponding to the baseline null model discussed above. For this purpose, the asymmetric summary measures of association based on the concept of Proportionate Reduction of Error (PRE) and already introduced above in Section 3.2.2, for example, the Goodman-Kruskal λ and τ , would appear as good candidates for evaluating prediction accuracy, as their premises suit the evaluation task at hand. Prediction and classification can be considered a one-way

relationship between the original data classes (as the independent variable) and the predicted, classified data classes (as the dependent variable), mediated by the explanatory variables included in the model, where the perfect relationship with all the instances on the diagonal would correspond with perfect prediction accuracy. Furthermore, in order to be of any actual worth we can rightly expect that the prediction or classification process on the basis of the models should exceed some baselines or thresholds, the levels of which correspond to the null relationships (Cohen et al. 2003: 516-519).

The problem with the original versions of these asymmetric association measures is that they do not distinguish between overall correct and incorrect classification; a perfect positive relationship receives the same association value as a perfect negative relationship (Menard 1995: 24-28). Fortunately, this can be remedied by slight adjustments to the formulas, where we compare prediction/classification errors *with* the model, ϵ_{model} , to the baseline level of prediction/classification errors *without* the model, $\epsilon_{\text{baseline}}$, according to formula 3.52 (Menard 1995: 28-30). The formula for the error with the model remains the same, irrespective of whether we are evaluating prediction or classification accuracy, presented in 3.49, but the errors without the model vary according to the intended objective, presented in 3.50 and 3.51. Subsequently, the measure for the *proportionate reduction of prediction error* is presented in 3.53, and, being analogous to the Goodman-Kruskal λ , it is designated as $\lambda_{\text{prediction}}$. This measure may maximally range between $[1-K, 1]$, with positive values indicating a better than baseline prediction, and negative values a worse performance. Similarly, the measure for *proportionate reduction of classification error* is presented in 3.54, and, being analogous with the Goodman-Kruskal τ , it is likewise designated as $\tau_{\text{classification}}$. This measure may range between $1-[K^2/(2\cdot K-2)]$ and 1.0, with positive values indicating a better than baseline classification, and negative values a worse performance; when the marginal (overall) distributions are unequal, as is the case here, the maximum value is less than one. Here, as well as with the original Goodman-Kruskal association measures, one should note that their ranges are not fixed, but will vary in accordance with the marginal distributions, which in this study are the original and predicted overall frequencies of the four lexemes.

$$(3.49) \epsilon_{\text{model}} = N - \sum_{k=1 \dots K} n_{k,k} = N - \sum \text{diag}(n), \text{ where } n \text{ is the prediction/classification matrix}$$

$$(3.50) \epsilon_{\text{baseline, prediction}} = N - \max(R_k), \text{ with } R_k = \sum_{i=1 \dots K} n_{k,i} \text{ for each row } k$$

$$(3.51) \epsilon_{\text{baseline, classification}} = \sum_{k=1 \dots K} \{R_k \cdot [(N-R_k)/N]\}, \text{ with } R_k = \sum_{i=1 \dots K} n_{k,i}$$

$$(3.52) \text{PRE} = (\epsilon_{\text{baseline}} - \epsilon_{\text{model}}) / \epsilon_{\text{baseline}}$$

$$(3.53) \lambda_{\text{prediction}} = 1 - \epsilon_{\text{model}} / \epsilon_{\text{baseline, prediction}}$$

$$(3.54) \tau_{\text{classification}} = 1 - \epsilon_{\text{model}} / \epsilon_{\text{baseline, classification}}$$

For these prediction and classification efficiency measures, one can even calculate the significance of the difference between the prediction errors with and without the model (Menard 1995: 30-31, 93, Note 10). I have implemented these in the *R* function calculating the presented prediction efficiency measures, but their specifics are beyond the primary scope of this dissertation.

Evaluating the robustness of the model effects with resampling schemes

There are various approaches with respect to what data the evaluation of prediction efficiency of the model will be calculated on. This evaluation of the prediction accuracy is often referred to as the *validation* of the model. In the first place, prediction can be undertaken on the original data (or various samples thereof), with which the model was fitted, which is called *internal validation* of the model. However, this obviously involves the risk of overestimating the accuracy of the model. To remedy this, in *external validation* one uses data which has not been originally used in training and fitting the models to evaluate the prediction efficiency of the fitted model. The simplest solution for external validation, known as *data-splitting* (Howell 2001: 90), is to set aside some portion of the data during the training and fitting stage, so that this leftover data can be considered “new” to the model at the testing stage; alternatively, one can acquire entirely new data for validation.

In this study, an obvious split would be to use the newspaper portion of the corpus for training the model, and the Internet newsgroup discussion portion for testing, or vice versa. However, this held-out or new data has to be sufficiently similar in its characteristics to the original data and of considerable size in absolute terms, in order to guarantee accurate evaluation and thus serve its purpose. For binary outcomes, the bare minimum is 100 cases for the less frequent outcome category, and even then reliability of the results is not guaranteed (Harrell 2001: 92). For polytomous outcomes one can expect the minimum per class to be at least as high, amounting to $4 \times 100 = 400$ instances of testing data (11.8% of all the 3404 instances of data) in this study with four lexemes. This means that a considerable amount of relevant information (and associated collection and analysis work) would have to be kept outside the fitting process and could play no role in the actual description of the studied phenomenon. Furthermore, setting aside a portion of the original data and thus diminishing the size of the training data may lead to undesirable restrictions in variable selection. Then, the parameters of the model will reflect only the effects evident in the training data and will certainly miss those present only in the testing data.

Completely new data naturally involves additional work in both acquiring and preparing it for statistical analysis, which may be significant in magnitude or unreasonably difficult, if not impossible, to accomplish (consider, e.g., data of historical linguistic usage, where one simply has to make do with what has survived in some recorded form), and it is thus often not a practically feasible option. A major advantage of splitting data is that it allows for testing hypotheses based on the training data with the testing data, but the disadvantages are considerable (Harrell 2001: 92-93). An alternative approach to external validation, noted earlier in Section 3.2.1, is not to gather more data, from similar sources with the same methods, but instead to pursue the same research question in an attempt to replicate the results with a different type of evidence, and thus also different methods, as suggested, for example, by Rosnow and Rosenthal (1989) and Moran (2003).

Resampling schemes are a remedy to the disadvantages of data-splitting, and they capitalize in an increasing manner on the capability of modern computers to sample

and analyze large data sets repeatedly in a (relatively) short time. The basic idea of resampling is to repeat the data-splitting and sampling process, of which there are various schemes to be discussed below, a considerable (say 10, 20, 50 or 100) or even extremely high (1000–10000) number of times, each time first fitting the model(s) with a newly sampled training portion of the data and, then, validating the particular result with the testing portion of the data. The evaluation of the model's performance is, in the end, based on the distribution of all the measures calculated for each of the individual testing portions, and summarized, for example, as an average, standard deviation, and/or confidence interval of the measure(s) of interest. The purpose of the repeated resampling is to ensure that all the data is taken into account both in the training and fitting as well as the testing of the models. If some phenomenon is present in the data, it should be represented in at least one (and possibly more) of the training and testing samples, and will thus contribute to the distribution of the measure of interest.

However, the overall value and weight of such a summary measure describing the model and its performance is dependent on how general the phenomenon is, which is reflected in how broadly it is present throughout the data and the individual samples. Therefore, such resampling schemes can be used to evaluate, not only of the performance of the model in predicting outcomes in the testing portions of the data, but also the robustness of the model itself through the accumulating estimates of the parameters (i.e., coefficients) of the explanatory variables included in the model, based on the training portions of the data. Furthermore, and most importantly, the variability of a selected phenomenon is in resampling studied through the data sample *at hand*, rather than by making assumptions concerning its distribution in the overall population, represented by so-called parameters (e.g., average, variance, and standard deviation), and trying to infer these from the sample. Therefore, resampling schemes provide *non-parametric* estimates, of both the model's description of the data and the model's performance in prediction, which neither require nor make any assumptions regarding the underlying population in its entirety.

The oldest resampling scheme is the *jack-knife*, also known as *cross-validation* (Mooney and Duval 1993: 22-27). In the jack-knife procedure, all the available data is divided (possibly, but not necessarily, randomly) into some predetermined number g of mutually exhaustive portions (which are thus samples *without replacement* of the original data), whereafter each portion is in turn left aside as the testing portion and all the remaining $g-1$ portions are used for training; consequently, this training and evaluation process is repeated g times. If the portions are split randomly, the entire process can be repeated, say, 10 times. At its extreme, the data is divided into as many portions as it contains individual instances, i.e., $g=N$, which is known as *leave-one-out* cross-validation. However, research indicates that grouped cross-validation, with $g=10$ or 20 , produces more accurate results than the leave-one-out procedure (Harrell 2001: 93). Nevertheless, any version of the jack-knife procedure leaves a portion of the data, albeit relatively small, outside the fitting and training stage during each iteration round, and thus the procedure cannot validate the model fully with the entire data.

The *bootstrap* procedure, introduced by Efron (1979; see also, e.g., Mooney and Duval 1993: 9-15; Fox 1997: 493-511; Harrell 2001: 87-90), offers a solution to this disadvantage, and it appears to have become the predominant resampling scheme in

recent years. In the bootstrap, one repeatedly selects random samples (of the same size N as the original data), *with* replacement, from the original data sample, with which the model is then trained and fitted. Consequently, each sample may contain some of the original instances more than once and some instances might not appear in each sample. Each of these fitted models is thereafter always tested with the entire original data sample. This process is repeated a substantial number of times, ranging from 50 upwards, depending on how much one wants to capitalize on the key characteristic of the bootstrap, to be described below. After the iterations are completed, the distributions of the calculated values of interest describe directly the data sample at hand, and indirectly give an indication of the studied phenomenon in the underlying population.

The central feature of the bootstrap is that, due to the nature of the resampling procedure, the distribution of any measure or descriptor of interest calculated concerning the data sample is constituted exactly of the set of the individually obtained values. Therefore, one does not have to make any assumptions about the distribution since it is available in its entirety, and any descriptive parameters can be calculated directly on the basis of this distribution. Furthermore, through resampling the original sample, the intention is to replicate the results of repeated sampling of the *underlying population*, and thus to asymptotically approach a *direct* estimate of the variation and distribution of the variables of interest in the original population. This is in contrast to estimating the probability of some values of such variables calculated from the original sample, given assumptions about their distribution in the population (Mooney and Duval 1993: 9-15, 20-22). So, if the number of repeated iterations is sufficiently large, at least $n \geq 1000$, one can calculate for a measure of interest, with some critical P-level (α), the associated confidence interval $\{p_{low} = \alpha/2, p_{high} = 1 - \alpha/2\}$ by simply sorting the values calculated for each iteration round and picking out the two with the indexes corresponding to the integer portions of the two percentiles, $n \cdot p_{low}$ and $n \cdot p_{high}$, respectively (known as the *percentile* method, see Mooney and Duval 1993: 36-37; Fox 1997: 503). In fact, it has been observed that the improvements in the accuracy of measures estimated with the bootstrap are only slight when the number of iterations rounds is increased to over 1000 (Mooney and Duval 1993: 21).

For some statistical procedures such as the fitting of logistic regression models, as is the case in this study, the combined effect of their iterated calculations may still take exceedingly long, despite ever increasing computational efficiency. In such circumstances (cf. Mooney and Duval 1993: 37) one can make do with a smaller number of iterations, $50 \geq n \geq 200$, and assume that the resampled values are distributed approximately normally. Then, having calculated descriptive measures such as the mean and variance⁷⁰ of the values of interest, one can approximate the confidence intervals according to the normal distribution. However, this *normal approximation* method (Mooney and Duval 1993: 33-36, Fox 1997: 502) is not generally recommendable, as it fails to take full advantage of the inherent nonparametric nature of the bootstrap procedure. Nonetheless, it may be the most practical solution when validating more complex heuristics or in multiple outcome settings with large

⁷⁰ Normally, the bootstrap mean is the simple mean of the bootstrapped values, but in some cases one might rather prefer to use the *trimmed* mean in order to reduce the influence of the outlying values, which can be quite extreme. Then, the variance is computed normally against the mean, whichever way it has been calculated.

numbers of constituent binary logistic regression models, especially if parallel computation resources are unavailable.⁷¹ Nevertheless, one should bear in mind that the bootstrap is better suited for estimating parameter value ranges such as confidence intervals, rather than exact points such as means/averages, since in the latter case outlying, extreme values can distort the result (Mooney and Duval 1993: 60). Furthermore, one should note that the simple bootstrap estimates exhibit some positive bias in favor of the models, for which a range of corrective measures have been presented (Mooney and Duval 1993: 37-42; Harrell 2001: 94-96); however, they, too, fall outside the scope of this study.

In the resampling process, the simplest method is to repeatedly sample randomly, *with replacement*, from the entire data sample as such; instances which have been sampled during one iteration are not put aside, but may be resampled during both the same iteration round and the next one(s) to the extent as chance allows.⁷² However, if one suspects that the original data sample might be clustered in such a way that individual groups may have sufficiently influential tendencies separating each one from the rest, and, furthermore, if it is not feasible to include this grouping as an explanatory variable to the model, this potential cluster-specific bias can be reduced and its effect assessed by sampling (with replacement) *within the groups* (Hoffman et al. 2001). For instance, this is the case when the data has been acquired in clusters (Fox 1997: 510), something which one can consider applicable for a corpus constituted by a large set of individual texts or utterances (being in the case of newspaper articles or Internet newsgroup postings relatively short both in length and in the time required to originally produce them).

In practice, what this means is that sampling is stratified so that each training portion of the data contains only one instance per group/cluster, with each instance being randomly sampled from within each group (with replacement, entailing that all the instances in each group are again available for random sampling during the next iteration round). This is a feasible method as long as the groups/clusters are relatively small in comparison to the entire data sample, so that the resultant training portions remain sufficiently large. However, the number of iteration rounds necessary for stable estimates might grow as high as 10000–50000 as the number of clusters

⁷¹ When I tested the fitting of a single binary model pitting *ajatella* against the rest, with 25 explanatory variables, using the simple bootstrap procedure on the entire available data sample (with 3404 instances) for 1000 times (which is required by the percentile confidence interval method), this took 20 hours on MacBook with 2GB memory and a 1.83 GHz Intel Core Duo Processor (see Section 3.4.5 below). As each heuristic for polytomous regression requires the fitting of several binary logistic regression models for each round, the overall duration with *serial* computation would turn out to be prohibitively large for the comparison of the various heuristics, especially in the case of ensembles of nested dichotomies.

⁷² This intuitive approach to sample directly from the observed outcomes and the associated values of the explanatory variables implicitly treats the selection of the contextual variables in the model as random rather than fixed (Fox 1997: 505). This would seem to suit this type of linguistic setting where we can hardly consider the selected set of variables as exhaustively and comprehensively determined for good, since we cannot *a priori* rule out that there could be further contextual features not included in the model, which might be relevant to linguistic usage not represented in the data at hand. In this sampling *from observations*, we evaluate whether the explanatory variables in the model are significantly relevant or not. Instead of the observations, however, one could alternatively sample *from residuals*, in which case the model and its selection of explanatory variables is implicitly considered to be correct (Mooney and Duval 1993: 16-17; Fox 1997: 506). Then, one rather attempts to mutually balance the weights of the variables, without questioning their inclusion and relevance in the model.

increases (Hoffman et al. 2001: 1125).⁷³ In a linguistic study, such clusters could be individual speakers/writers, if the data sample is pooled from a large number of their utterances and texts, or it could be individual discourse passages or texts from which the sample corpus is compiled. In accordance with Bresnan et al. (2007), one could very well hypothesize that individual speakers/writers or individual fragments of discourse/text may exhibit preferences, the effects of which on the actual model one would like to assess with such sampling.⁷⁴

If the size of the groupings in relation to the entire original sample grows, the size of the stratified sample in the within-groups resampling scheme decreases prohibitively; an example of this in a linguistic context could be high-level classifications of a corpus such as text types, genre, mode (e.g., spoken vs. written) or medium (e.g., published vs. Internet). In such a case, one possible solution is to use *stratum/cluster-based resampling*, where one samples with replacement from each cluster/stratum individually as many instances as there originally are in each original cluster/stratum (Fox 1997: 510-511), which is an approach suggested by Gries (2007).⁷⁵ However, Hoffman et al. (2001) indicate that, of the various cluster-based schemes, only within-cluster resampling also remains valid when the cluster-related effect is real and *non-ignorable*. An alternative approach is to treat such a grouping as an explanatory variable incorporated in the model, which is what Bresnan et al. (2007) also did.

In this study, I will first analyze and evaluate the data sample using the simple bootstrap, without assuming writer bias, but I will follow this initial analysis by a second one in which the writers are treated as clusters and resampled accordingly. Furthermore, in a third analysis I will treat the *medium* (newspaper article vs. Internet discussion) as an explanatory variable. Together, these three types of analyses should shed light on the potential interaction of the strictly linguistic variables with the extralinguistic effects.

⁷³ Since the individual resamples and associated fits are independent of each other, this task is a perfect candidate for parallel computation. Especially at such higher magnitudes of iterations, since the overall duration in sequential calculation is a simple multiple of the time required for a singular fit of the model, dividing the task over multiple processors and thus fitting models concurrently can drastically reduce the required time. In fact, with the valuable assistance of Jarno Tuimala and Yrjö Leino at CSC – IT Center for Science <www.csc.fi>, I have implemented a simple script scheme by which this can be run on CSC’s parallel workstations. Consequently, computing a 10000-fold bootstrap with resampling from clusters by dividing the task into 400-fold iterations running parallel on 25 processors, the overall duration is reduced to 4% of a corresponding sequential computation. On CSC’s *Murska* parallel computer, an HP CP4000 BL ProLiant supercluster with 2176 compute cores, or 1088 dual-core 2.6 GHz AMD Opteron 64-bit processors, each 400-fold iteration, and thus the entire 10000-fold computation, of the full polytomous logistic regression model (with all four THINK lexemes) took approximately only 17-18 minutes.

⁷⁴ Bresnan et al. (2007) report that they use “*bootstrap sampling with replacement of entire clusters* [i.e. speakers]. ... The same speaker’s data can randomly occur many times in each copy. ...”. If I understand this to mean that the sampling process concerns clusters of several instances at a time, this could lead to variation in the overall resample size, but it is not stated what steps, if any, are then taken to make the size of the resample exactly equal to the original sample. Alternatively, this could be understood to refer to *cluster/stratum-based resampling* of the type used by Gries (2007), or even *within-cluster resampling* of the type suggested by Hoffman et al. (2001), but I cannot discern this on the basis of what is explicitly put forth.

⁷⁵ In fact, Gries (2007) compares the *exhaustive permutation of clusters* with bootstrapping using cluster-based resampling, and argues in favor of the latter method, because it is applicable at any level of granularity but does not run the risk of becoming computationally unfeasible as the number of partitions grows rapidly (at the exponential rate of $2^n - 1$).

3.4.5 A detailed example of a (binary) logistic regression model

For the purpose of illustrating in detail how logistic regression works and what results it produces, and thus what the aggregated (and somewhat simplified and summarized) results of the various polytomous heuristics are based on, I will present one binary logistic model, namely, contrasting the occurrence of *ajatella* against the other three, less frequent THINK lexemes. As explanatory variables at this time, I have selected all those which have been discussed explicitly above in the presentation of univariate and bivariate methods in Sections 3.2 and 3.2, and which have a sufficient overall frequency and occurrences with more than only one of the studied lexemes (see Tables 3.14, 3.21, and 3.29 above). These are 1) the six person/number features, 2) the two semantic types of the AGENT, and 3) the 17 semantic and structural types of the PATIENT, adding up to 25 variables in all. This model has been fitted with the entire data sample (with 3404 instances) by the standard `glm` function available in *R*, and the estimates of the coefficients and their significances are presented in Table 3.43. The prediction efficiency of the resultant model has in this case been evaluated with the original training data.

```
glm(formula = ajatella ~ Z_SG1 + Z_SG2 + Z_SG3 + Z_PL1 + Z_PL2
+ Z_PL3 + SX_AGE.SEM_INDIVIDUAL + SX_AGE.SEM_GROUP +
SX_PAT.SEM_INDIVIDUAL + SX_PAT.SEM_GROUP + SX_PAT.SEM_NOTION +
SX_PAT.SEM_ATTRIBUTE + SX_PAT.SEM_STATE + SX_PAT.SEM_TIME +
SX_PAT.SEM_ACTIVITY + SX_PAT.SEM_EVENT +
SX_PAT.SEM_COMMUNICATION + SX_PAT.SEM_COGNITION +
SX_PAT.SEM_LOCATION + SX_PAT.SEM_ARTIFACT +
SX_PAT.INDIRECT_QUESTION + SX_PAT.DIRECT_QUOTE +
SX_PAT.INFINITIVE + SX_PAT.PARTICIPLE + SX_LX_että_CS.SX_PAT,
family = binomial, data = THINK.data)
```

Table 3.43. Parameter values and other associated statistics of the fitted binary logistic regression model contrasting *ajatella* against the three other THINK lexemes, with person/number, semantic types of AGENT, and semantic and structural types of PATIENT as explanatory variables, adapted from `glm(...)` output in *R*. Significant (with $P < 0.05$) odds-ratios of variables in boldface; significance codes: ‘***’ ~ $P < 0.001$, ‘**’ ~ $P < 0.01$, ‘*’ ~ $P < 0.05$, ‘.’ ~ $P < 0.1$, ‘-’ ~ $P > 0.1$.

Explanatory variables/ Coefficients	Odds	Log- odds	Std. Error (ASE)	z-value	P(> z)	Sign. code
(Intercept)	2.071	0.728	0.010	7.307	$2.725e^{-13}$	***
Z_SG1	1.965	0.676	0.180	3.759	0.0002	***
Z_SG2	1.286	0.251	0.198	1.269	0.2045	–
Z_SG3	1.024	0.023	0.144	0.162	0.8710	–
Z_PL1	3.716	1.313	0.569	2.306	0.0211	*
Z_PL2	0.584	-0.538	0.353	-1.525	0.1272	–
Z_PL3	1.834	0.607	0.208	2.922	0.0035	**
SX_AGE.SEM_INDIV...	0.855	-0.156	0.105	-1.491	0.1360	–
SX_AGE.SEM_GROUP	0.232	-1.459	0.224	-6.514	$7.337e^{-11}$	***
SX_PAT.SEM_INDIV...	1.691	0.526	0.265	1.981	0.0476	*
SX_PAT.SEM_GROUP	5.477	1.701	0.608	2.796	0.0052	**
SX_PAT.SEM_NOTION	0.179	-1.720	0.123	-14.039	$<2e^{-16}$	***
SX_PAT.SEM_ATTR...	0.200	-1.608	0.287	-5.597	$2.186e^{-08}$	***
SX_PAT.SEM_STATE	0.490	-0.714	0.356	-2.006	0.0449	*
SX_PAT.SEM_TIME	0.731	-0.313	0.349	-0.896	0.3703	–
SX_PAT.SEM_ACT...	0.120	-2.119	0.142	-14.868	$<2e^{-16}$	***
SX_PAT.SEM_EVENT	1.154	0.142	0.414	0.345	0.7300	–
SX_PAT.SEM_COMM...	0.084	-2.480	0.448	-5.533	$3.150e^{-08}$	***
SX_PAT.SEM_COGN...	0.418	-0.872	0.482	-1.808	0.07058	.
SX_PAT.SEM_LOC...	1.346	0.297	0.541	0.549	0.5831	–
SX_PAT.SEM_ARTIFACT	1.436	0.362	0.584	0.620	0.5354	–
SX_PAT.INDIRECT Q...	0.049	-3.015	0.185	-16.319	$<2e^{-16}$	***
SX_PAT.DIR... QUOTE	0.015	-4.170	0.601	-6.939	$3.942e^{-12}$	***
SX_PAT.INFINITIVE	4.904	1.590	0.543	2.930	0.003389	**
SX_PAT.PARTICIPLE	4.474	1.498	0.371	4.033	$5.498e^{-05}$	***
SX_LX että CS.SX PAT	1.924	0.655	0.147	4.462	$8.131e^{-06}$	***

The null deviance D_{null} , based on only the overall relative proportion of the lexeme *ajatella*, is 4667.0, and the model deviance D_{model} , remaining after the explanatory variables are taken into consideration, is 3347.3. Thus, the relative decrease in deviance, reflecting the fit of the model with the data, is $R_L^2 = 1 - (3347.3/4667.0) = 0.283$. This is not a bad fit at all, considering that the variables included in this example represent only a subset of all the potential ones, though they probably do include the most important ones in the case of the lexemes in question. With respect to prediction efficiency, based on the prediction table presented in Table 3.44, the overall recall rate was 77.2%, while the measures assessing the reduction of error are $\lambda_{prediction} = 0.480$ and $\tau_{classification} = 0.537$, which are also quite good results. Lexeme-wise, the *Recall* for *ajatella* was 78.3% and the *Precision* 72.1%, whereas for the other lexemes as a group the *Recall* was 76.4% and the *Precision* 81.8%.

Table 3.44. Prediction table of *ajatella* vs. the rest resulting from the selected explanatory variables.

THINK.one_vs_rest.A_vs_other.Z_PERSON_NUMBER.SX_AGE_PAT\${test.guess.mean, test.guess.rel, test.lx, guess.lx, success.lx}

Observed/Predicted	<i>ajatella</i>	Other	Σ(Observed Lexeme)
<i>ajatella</i>	1170 (78.4%)	322 (21.6%)	1492
Other	453 (23.7%)	1459 (76.3%)	1912
Σ(Predicted Lexeme)	1620	1784	3404

Turning to the fitted model, as many as 16 of the altogether 25 coefficients (in addition to the intercept) – corresponding to the selected variables – were assessed as significant (on the basis of comparing the log-odds values with their asymptotic standard errors in the data). Of these, 8 had positive and 8 had negative log-odds values, consisting on the one hand of the strongest odds for the occurrence of *ajatella* in association with a human GROUP as PATIENT (5.477), and, on the other hand, the greatest odds (0.015) against its occurrence with a DIRECT QUOTE as a PATIENT. With regard to specific feature groups, the FIRST PERSON PLURAL (3.716), FIRST PERSON SINGULAR (1.965), and THIRD PERSON PLURAL (1.834) features, in descending order, were associated with *ajatella*, with the other person/number features remaining neutral. With respect to the two types of AGENT under scrutiny, human GROUPS decrease the odds (0.232 ~ 1:4.3) of *ajatella* occurring, whereas human INDIVIDUALS are not a significantly distinctive feature as an AGENT. For the different semantic and structural types of PATIENT, INFINITIVES, PARTICIPLES, and the *että*-clause (‘that’) in addition to human referents, whether INDIVIDUALS or GROUPS, show positive odds for *ajatella*, while abstract NOTIONS, ATTRIBUTES, and STATES, ACTIVITIES and acts/forms of COMMUNICATION, as well as both INDIRECT QUESTIONS and DIRECT QUOTES decrease the odds, to differing degrees. Nevertheless, we should remember that the individual semantic and structural types within each feature group studied here are mutually exclusive, and the results are in effect based on (maximally) feature trios with one feature each from of 1) person/number, 2) AGENT type, and 3) PATIENT type.

We can now make a preliminary comparison between these multivariate results and those gained with the univariate analyses presented earlier above, focusing on the relationship of the selected feature variables and the occurrence of the lexeme *ajatella* (and disregarding the three other THINK lexemes until the full-scale analysis to follow later on), laid out in Table 3.45. What becomes clear is that there is a definite correspondence between the two levels of analysis, though this relationship is not categorical. In the case of 13 variables, positive as well as negative associations assessed as significant in the multivariate analysis are matched by similar preferences in the earlier univariate analyses. Nevertheless, a substantial proportion of features considered significant in the univariate analyses do not turn out to be so in the multivariate analysis, when these variables are considered in relation to each other in their entirety (THIRD PERSON SINGULAR, SECOND PERSON PLURAL, human INDIVIDUALS as AGENT, and TIME, EVENT, LOCATION, and ARTIFACT as PATIENT), and the other way around (FIRST PERSON PLURAL, and STATE and ATTRIBUTE as PATIENT). Even so, no associations are observed to have become reversed between the univariate and multivariate analyses. Furthermore, the magnitudes of the associations are not always of similar strength, for instance, in the case of *että*-clauses, though INDIRECT QUESTIONS show that this divergence in results is not categorical.

Table 3.45. Comparison of the univariate results, based on standardized Pearson residuals (e_{ij}) of the distribution of the selected features among the studied lexemes (derived from Tables 3.19, 3.26, and `multiple.feature.distribution` (THINK.data, think.lex, SX_PAT.classes) `$residual.pearson.std.sig`), and the multivariate results based on the logistic regression model of these same features (derived from Table 3.43), with respect to the occurrence of *ajatella* against the rest. Significant values are set in boldface, with significant positive association with *ajatella* indicated by '+', a negative positive association with '-', and a nonsignificant result with '0'.

Feature/Measure (<i>ajatella</i>)	Univariate result	Stand. Pearson residual	Odds	Multivariate result
Z_SG1	+	+7.636	1.965	+
Z_SG2	+	+2.073	1.286	0
Z_SG3	-	-9.072	1.024	0
Z_PL1	0	+1.815	3.716	+
Z_PL2	-	-2.011	0.584	0
Z_PL3	+	+2.328	1.834	+
SX_AGE.SEM_INDIVIDUAL	+	+9.811	0.855	0
SX_AGE.SEM_GROUP	-	-9.811	0.232	-
SX_PAT.SEM_INDIVIDUAL	+	7.076	1.691	+
SX_PAT.SEM_GROUP	+	6.049	5.477	+
SX_PAT.SEM_NOTION	-	-6.003	0.179	-
SX_PAT.SEM_ATTRIBUTE	0	-1.489	0.200	-
SX_PAT.SEM_STATE	0	1.136	0.490	-
SX_PAT.SEM_TIME	+	2.573	0.731	0
SX_PAT.SEM_ACTIVITY	-	-9.520	0.120	-
SX_PAT.SEM_EVENT	+	3.795	1.154	0
SX_PAT.SEM_COMMUNICATION	-	-2.892	0.084	-
SX_PAT.SEM_COGNITION	0	0.801	0.418	0
SX_PAT.SEM_LOCATION	+	3.273	1.346	0
SX_PAT.SEM_ARTIFACT	+	3.318	1.436	0
SX_PAT.INDIRECT_QUESTION	-	-12.895	0.049	-
SX_PAT.DIRECT_QUOTE	-	-7.733	0.015	-
SX_PAT.INFINITIVE	+	7.518	4.904	+
SX_PAT.PARTICIPLE	+	9.563	4.474	+
SX_LX_että_CS.SX_PAT	+	20.221	1.924	+

Next, we can use the bootstrap as an alternative way to assess the significance of the effects in the data, represented by the coefficients β_m (i.e., log-odds) associated with the explanatory variables X_m . Keeping the critical P-value as $\alpha=.05$, we can construct the corresponding 95% percent confidence interval with the percentile method by fitting the model repeatedly according to the simple bootstrap sampling, making 1000 iteration rounds in order to enable us to use the percentile method to produce the low and high estimate values.⁷⁶ The results presented in Table 3.46 show that the 95% confidence intervals are quite broad, and in a few cases (human INDIVIDUALS and STATES as PATIENTS) the effects are no longer significant as the intervals bridge both sides of the odds-ratio $\exp(\beta_m) = 1$ (i.e., the null odds of 1:1, corresponding to the log-

⁷⁶ This is the procedure mentioned above which took 20 hours to complete on a current laptop computer using serial computation. Thus, one is not tempted to make a habit of resorting to it just to "check things out". However, having access to a parallel computer such as CSC's *murska*, the effective duration can reduced to as little as 3-4 minutes using 20 concurrently fitted 50-fold partitions.

odds $\beta_m=0$). Furthermore, in a few cases, especially with DIRECT QUOTES as PATIENTS, the upper end of the confidence interval for the odds-ratio is absurdly high (e.g., $23983545 \approx 2.4e^{10}$, though the corresponding logs-odds are somewhat more reasonable at 17.0).

If such values are merely chance quirks, they should get eliminated by the percentile method. So, this is indicative of some of the difficulties in fitting this model, which may possibly result from some close to exact correlation among some of the variables, an aspect which was not accounted for in this example case, in addition to extremely skewed distributions of the features in question due to the random sampling process. Finally, we can also calculate the 95% Confidence Intervals for other statistics evaluating the fit with respect to the entire data and the prediction efficiency of the model, which are $R_L^2=(0.255, 0.280)$, $\lambda_{prediction}=(0.474, 0.485)$, $\tau_{classification}=(0.532, 0.542)$, and overall $Recall=(76.94\ 77.41\%)$. Taking the lexeme-wise perspective, the 95% Confidence Intervals are, in the case of *ajatella*, (0.763, 0.794) for *Recall* and (0.716, 0.728) for *Precision*, while for the other THINK lexemes, when lumped together, the value ranges are (0.755, 0.777) for *Recall* and (0.807, 0.825) for *Precision*.

```
THINK.one_vs_rest.A_vs_other.Z_PERSON_NUMBER.SX_AGE_PAT.1000
<-
polytomous.logistic.regression(data.internal=THINK.A_vs_other.
data,,fn="Z_SG1 + Z_SG2 + Z_SG3 + Z_PL1 + Z_PL2 + Z_PL3 +
SX_AGE.SEM_INDIVIDUAL + SX_AGE.SEM_GROUP +
SX_PAT.SEM_INDIVIDUAL + SX_PAT.SEM_GROUP + SX_PAT.SEM_NOTION +
SX_PAT.SEM_ATTRIBUTE + SX_PAT.SEM_STATE + SX_PAT.SEM_TIME +
SX_PAT.SEM_ACTIVITY + SX_PAT.SEM_EVENT +
SX_PAT.SEM_COMMUNICATION + SX_PAT.SEM_COGNITION +
SX_PAT.SEM_LOCATION + SX_PAT.SEM_ARTIFACT +
SX_PAT.INDIRECT_QUESTION + SX_PAT.DIRECT_QUOTE +
SX_PAT.INFINITIVE + SX_PAT.PARTICIPLE + SX_LX_että_CS.SX_PAT",
lex=c("ajatella","other"),, classifier="one.vs.rest",
validation="internal.boot.simple", iter=1000,
ci.method="percentile",trim=.5)
```

Table 3.46. Confidence intervals ($CI=95\% \Leftrightarrow \alpha=0.05$; $CI: \alpha/2 < \exp(\beta_m) < 1-\alpha/2$), calculated with the percentile method using a simple bootstrap repeated 1000 times, of coefficients of the fitted binary logistic regression model contrasting *ajatella* against the three other THINK lexemes, with person/number, semantic types of AGENT, and semantic and structural types of PATIENT as explanatory variables; significant ranges of odds-ratios (with entire $CI < 1$ or $CI > 1$) of variables in boldface; results differing from the original single-round fit with the entire data in italic and with thicker border-lines.

THINK.one_vs_rest.A_vs_other.Z_PERSON_NUMBER.SX_AGE_PAT.1000\$odds.range

Feature/Lexeme	<i>ajatella</i>	<i>other</i>
(Intercept)	1.708<..2.545	0.391<..0.586
Z SG1	1.405<..2.846	0.351<..0.708
Z SG2	(0.785<.. 2.045)	(0.484<.. 1.274)
Z SG3	(0.7542<.. 1.392)	(0.719<.. 1.326)
Z PL1	1.031<..14	0.0703<..0.969
Z PL2	(0.245<.. 1.157)	(0.864<.. 4.076)
Z PL3	1.304<..2.544	0.393<..0.764
SX_AGE.SEM_INDIVIDUAL	<i>(0.694<..1.058)</i>	<i>(0.944<..1.441)</i>
SX_AGE.SEM_GROUP	0.145<..0.345	2.899<..6.832
SX_PAT.SEM_INDIVIDUAL	1.001<..3.033	0.330<..0.999
SX_PAT.SEM_GROUP	2.078<..34	0.0293<..0.481
SX_PAT.SEM_NOTION	0.140<..0.229	4.35<..7.161
SX_PAT.SEM_ATTRIBUTE	0.1<..0.330	3.036<..9.918
SX_PAT.SEM_STATE	<i>(0.243<..1.035)</i>	<i>(0.944<..4.059)</i>
SX_PAT.SEM_TIME	(0.372<.. 1.599)	(0.619<.. 2.687)
SX_PAT.SEM_ACTIVITY	0.0862<..0.159	6.272<..11
SX_PAT.SEM_EVENT	(0.530<.. 3.078)	(0.324<.. 1.886)
SX_PAT.SEM_COMMUNICATION	0.0262<..0.176	5.668<..37
SX_PAT.SEM_COGNITION	(0.125<.. 1.401)	(0.707<.. 8.031)
SX_PAT.SEM_LOCATION	(0.488<.. 6.007)	(0.166<.. 1.974)
SX_PAT.SEM_ARTIFACT	(0.462<.. 6.789)	(0.146<.. 2.151)
SX_PAT.INDIRECT_QUESTION	0.0303<..0.0693	14<..33
SX_PAT.DIRECT_QUOTE	0<..0.0371	27<..2.4e⁷
SX_PAT.INFINITIVE	2.289<..22	0.0425<..0.435
SX_PAT.PARTICIPLE	2.462<..11	0.0901<..0.403
SX_LX_että_CS.SX_PAT	1.439<..2.645	0.378<..0.695

We can now move on to evaluate whether writer/speaker-specific preferences have an influence on the results. This is done by repeating the bootstrap procedure for estimating confidence intervals, but this time, using the within-cluster resampling scheme with each writer/speaker (amounting to 571 in all) in the data interpreted as a single cluster, in the spirit of Bresnan et al. (2007). As it is recommendable to use more iterations for this scheme than what is required for the simple bootstrap, especially when the number of clusters is high, I will use 10000 repetitions here.⁷⁷ With this adjustment, the results presented in Table 3.47 below show that all the person/number features which showed significant association for *ajatella* (FIRST PERSON SINGULAR, FIRST PERSON PLURAL, and THIRD PERSON PLURAL) appear instead

⁷⁷ The overall duration of this scheme with sampling from the writers as clusters appears to have a similar ratio to the number of iterations as the simple bootstrap. Thus, computing 10000-fold repetitions serially on a standard desktop computer would require several days, so a parallel solution becomes a practical necessity. On the *murska* supercluster at CSC, partitioning the entire task into 40 parallel 250-fold iterations took only roughly 12-13 minutes to calculate.

subject to writer preferences. With respect to AGENT types, human GROUPS remain a significant feature associated with the lexemes other than *ajatella*, whereas with PATIENT types, human INDIVIDUALS, STATES, as was the case in the simple bootstrap, and now also INFINITIVES, too, are not significant in this writer-cluster bootstrap scheme, when compared to the basic model fit once with the entire data. Nevertheless, 10 out of 16 variables judged significant in the simple fit remain so even in this analysis, indicating that these variables represent robust effects, the identification of which is exactly the purpose of the writer-cluster bootstrapping scheme.

Furthermore, the confidence intervals for other statistics evaluating the prediction efficiency of the model are $\lambda_{prediction}=(0.440, 0.483)$, and $\tau_{classification}=(0.501, 0.540)$. The confidence interval for overall *Recall* is 75.44–77.35%, while the lexeme-specific *Recall* is (0.704, 0.727) and the *Precision* is (0.732, 0.804) for *ajatella*, whereas the corresponding values for the other THINK lexemes when lumped together are (0.739, 0.782) and (0.787, 0.829), respectively. All of these values are not in any practically significant extent different from than the ones derived with the simple bootstrap.

However, R_L^2 , which assesses the fit of the model with new, unseen data, provides dismal results this time, with a confidence interval of (-0.207, 0.226). It appears that as the fit with the considerably smaller clustered training data ($n=571$) can become quite high with $R_L^2=(0.276-0.394)$, the estimated odds then become stronger due to overfitting, which results in the range of the estimated probabilities becoming more extreme. This punishes the fit with the testing data, since it is in particular the lower probabilities (incorrectly) estimated for individual actual occurrences which most increase model deviance, which is reflected in the R_L^2 measure. In fact, if roughly one-third (i.e., 1000/3404) of the observations in this case receive a small probability estimate of $P \leq 0.097 = \exp(D_{null}/-2 \cdot 1000)$, the model deviance D_{model} is already worse than that for the null model, regardless of how good the probability estimates for the remaining two-thirds are, even if all these others were the maximum possible $P=1.0$. The same is also the case if only as few as 100 (2.9%) observations receive extremely bad probability estimates with $P \approx 0 (=7.34e^{-11} = \exp[D_{null}/-2 \cdot 100])$. Nevertheless, in a similar manner, one could just as well assess the influence of other extralinguistic factors manifested as small clusters, such as coherence within individual fragments of text or discourse, by resampling in such a case from each text/passage as a cluster.

```
polytomous.logistic.regression(data.internal=THINK.A_vs_other.
data,,fn="Z_SG1 + Z_SG2 + Z_SG3 + Z_PL1 + Z_PL2 + Z_PL3 +
SX_AGE.SEM_INDIVIDUAL + SX_AGE.SEM_GROUP +
SX_PAT.SEM_INDIVIDUAL + SX_PAT.SEM_GROUP + SX_PAT.SEM_NOTION +
SX_PAT.SEM_ATTRIBUTE + SX_PAT.SEM_STATE + SX_PAT.SEM_TIME +
SX_PAT.SEM_ACTIVITY + SX_PAT.SEM_EVENT +
SX_PAT.SEM_COMMUNICATION + SX_PAT.SEM_COGNITION +
SX_PAT.SEM_LOCATION + SX_PAT.SEM_ARTIFACT +
SX_PAT.INDIRECT_QUESTION + SX_PAT.DIRECT_QUOTE +
SX_PAT.INFINITIVE + SX_PAT.PARTICIPLE + SX_LX_ettà_CS.SX_PAT",
lex=c("ajatella","other"), freq, classifier="one.vs.rest",
validation="internal.cluster.speaker", iter=10000,
ci.method="percentile", trim=.5)
```

Table 3.47. Confidence intervals ($CI=95\% \Leftrightarrow \alpha=0.05$; $CI: \alpha/2 < \exp(\beta_m) < 1-\alpha/2$), calculated with the percentile method using simple bootstrap repeated 10000 times resampling from clusters, of coefficients of the fitted binary logistic regression model contrasting *ajatella* against the three other THINK lexemes, with person/number, semantic types of AGENT, and semantic and structural types of PATIENT as explanatory variables; significant ranges of odds-ratios (with entire $CI < 1$ or $CI > 1$) of variables in boldface; results differing from the original single-round fit with the entire data in italic and with thicker border-lines.

THINK.one_vs_rest.A_vs_other.Z_PERSON_NUMBER.SX_AGE_PAT.10000_speaker_cluster\$odds.range

Feature/Lexeme	ajatella	other
(Intercept)	1.2<..2.92	0.343<..0.835
Z SG1	(0.961<.. 6.42)	(0.156<.. 1.035)
Z SG2	(0.519<.. 6.19)	(0.161<.. 1.93)
Z SG3	(0.458<.. 1.88)	(0.532<.. 2.19)
Z PL1	(0.416<.. 2.5e⁸)	(0<.. 2.4)
Z PL2	(0<.. 6.02)	(0.166<.. 7.9e⁷)
Z PL3	(0.582<.. 3.22)	(0.31<.. 1.72)
SX AGE.SEM INDIVIDUAL	(0.614<.. 1.61)	(0.623<.. 1.63)
SX AGE.SEM GROUP	0.0521<..0.496	2.01<..19
SX PAT.SEM INDIVIDUAL	(0.681<.. 2.6e⁷)	(0<.. 1.44)
SX PAT.SEM GROUP	1.41<..5.9e¹²	0<..0.7
SX PAT.SEM NOTION	0.105<..0.312	3.21<..9.46
SX PAT.SEM ATTRIBUTE	0<..0.437	2.29<..5.4e⁷
SX PAT.SEM STATE	(0.158<.. 2.1e⁷)	(0<.. 6.34)
SX PAT.SEM TIME	(0.184<.. 5.5e⁷)	(0<.. 5.43)
SX PAT.SEM ACTIVITY	0.0596<..0.219	4.57<..17
SX PAT.SEM EVENT	(0.268<.. 3.3e⁷)	(0<.. 3.73)
SX PAT.SEM COMMUNICATION	0<..0.425	2.35<..1.5e⁸
SX PAT.SEM COGNITION	(0<.. 2.1e⁷)	(0<.. 1.5e⁸)
SX PAT.SEM LOCATION	(0.386<.. 3.3e⁷)	(0<.. 2.58)
SX PAT.SEM ARTIFACT	(0.162<.. 3.7e⁷)	(0<.. 6.17)
SX PAT.INDIRECT QUESTION	0.012<..0.0858	12<..84
SX PAT.DIRECT QUOTE	0<..0.0942	11<..1.7e⁸
SX PAT.INFINITIVE	(0.697<.. 4.5e⁷)	(0<.. 1.41)
SX PAT.PARTICIPLE	1.82<..4.6e⁷	0<..0.549
SX_LX_että_CS.SX_PAT	1.042<..4.47	0.224<..0.96

Finally, we can assess whether the *medium* of language usage has any significant effect, on top of the already selected explanatory variables, on the selection of *ajatella* in contrast to the other three THINK lexemes. Because I will in general not be including interaction effects in the models in this study due to the limiting sample size, I will only study the impact of including one more variable, representing the linguistic *medium*, on the fit and prediction efficiency of the model with the data.⁷⁸ The medium variable, denoted by the label Z_EXTRA_SRC_hs95, will be TRUE, if

⁷⁸ Because this particular dichotomous setting, pitting *ajatella* against the rest, with a higher limiting sampling size ($m=1492/10 \approx 150$) in comparison to the entire polytomous setting, nevertheless allows for a higher number of explanatory variables to be included in a model, for curiosity's sake I tried out a model with the medium variable in interaction with all the other variables, presented in Appendix M. While many of explanatory variables are not swayed by the medium, the results indicate that some others are, which is not, in the end, that surprising. However, the small frequency of *harkita* does not allow for studying interactions in the entire polytomous setting at hand.

the instance in question appears in the newspaper portion of the data, whereas the value will be FALSE, if the instance is to be found in the Internet newsgroup discussion portion. The statistics evaluating the fit and the prediction efficiency of the model are $R_L^2=0.292$, $\lambda_{prediction}=0.489$, and $\tau_{classification}=0.545$. The overall *Recall* is 77.61%, while the lexeme-specific *Recall* is 79.16% and the *Precision* is 72.37% for *ajatella*, whereas the corresponding values for the other THINK lexemes when lumped together are 76.41% and 82.45%, respectively. All of these values are higher than the ones for the model without a variable for the medium, as could be expected from adding an explanatory variable, but the increase is only slight.

The impact of the added variable on the relative weights of the other variables in the model is considerably greater (Table 3.48). Not only is the Z_EXTRA_SRC_hs95 feature significant in itself, with the corresponding odds (0.54561262) being in favor of the other lexemes (remember that these include the bookish *pohtia* and *harkita* in addition to the more common *miettiä*), but the number of other features with significant odds increases from 16 to 18, in comparison to the simple fit of the model with the entire data. These newly significant features are the SECOND PERSON PLURAL and human INDIVIDUALS as AGENT, both favoring the other three lexemes instead of *ajatella*, with the odds 0.502 and 0.767, respectively. These new developments, however, following from the inclusion of usage medium as also a variable, are not reversals, since in the previous assessments their effects have been regarded as insignificant. Nevertheless, this new model with a variable for usage medium could be further subjected to the same validation processes with the different bootstrap schemes as demonstrated above, and I have little doubt that this would not result in the decrease of the number of explanatory variables with a significant (robust) effect, more or less along the trend which was observed in the case of the slightly simpler model before.

```
polytomous.logistic.regression(data.internal=THINK.A_vs_other.
data,,fn="Z_SG1 + Z_SG2 + Z_SG3 + Z_PL1 + Z_PL2 + Z_PL3 +
SX_AGE.SEM_INDIVIDUAL + SX_AGE.SEM_GROUP +
SX_PAT.SEM_INDIVIDUAL + SX_PAT.SEM_GROUP + SX_PAT.SEM_NOTION +
SX_PAT.SEM_ATTRIBUTE + SX_PAT.SEM_STATE + SX_PAT.SEM_TIME +
SX_PAT.SEM_ACTIVITY + SX_PAT.SEM_EVENT +
SX_PAT.SEM_COMMUNICATION + SX_PAT.SEM_COGNITION +
SX_PAT.SEM_LOCATION + SX_PAT.SEM_ARTIFACT +
SX_PAT.INDIRECT_QUESTION + SX_PAT.DIRECT_QUOTE +
SX_PAT.INFINITIVE + SX_PAT.PARTICIPLE + SX_LX_että_CS.SX_PAT +
Z_EXTRA_SRC_hs95", lex=c("ajatella","other"), freq,
classifier="one.vs.rest", validation="internal.simple",
iter=1, ci.method="normal",trim=0)
```

Table 3.48. Coefficients and associated P-values of the fitted binary logistic regression model contrasting *ajatella* against the three other THINK lexemes, with *medium* in addition to person/number, semantic types of AGENT, and semantic and structural types of PATIENT as explanatory variables; significant values in boldface; results differing from the original single-round fit with the entire data in italics and with thicker border-lines.

Feature/Lexeme	<i>ajatella</i>	P-value
(Intercept)	2.733	0.0
Z EXTRA_SRC_hs95	0.546	0.0
Z SG1	1.960	0.000173
Z SG2	1.097	0.641
Z SG3	1.010	0.518
Z PL1	4.672	0.00774
Z PL2	<i>0.502</i>	0.0493
Z PL3	2.094	0.000528
SX AGE.SEM INDIVIDUAL	<i>0.767</i>	0.0134
SX AGE.SEM GROUP	0.250	0.0
SX PAT.SEM INDIVIDUAL	1.847	0.0224
SX PAT.SEM GROUP	6.523	0.00217
SX PAT.SEM NOTION	0.197	0.0
SX PAT.SEM ATTRIBUTE	0.228	0.0
SX PAT.SEM STATE	0.496	0.0516
SX PAT.SEM TIME	0.926	0.827
SX PAT.SEM ACTIVITY	0.150	0.0
SX PAT.SEM EVENT	1.520	0.316
SX PAT.SEM COMMUNICATION	0.0849	0.0
SX PAT.SEM COGNITION	0.407	0.0643
SX PAT.SEM LOCATION	1.720	0.321
SX PAT.SEM ARTIFACT	2.045	0.224
SX PAT.INDIRECT QUESTION	0.0543	0.0
SX PAT.DIRECT QUOTE	0.0222	0.0
SX PAT.INFINITIVE	5.422	0.00190
SX PAT.PARTICIPLE	4.485	0.000057
SX LX että CS.SX PAT	2.073	0.000001

To sum up, we can now compare the results of the various fitting and sampling schemes until now, presented in Table 3.49. Of the 25 originally selected explanatory variables, 10 remained significant throughout all the analyses, which suggests that the features in question most probably represent robust effects. Of these, three had odds-ratios in favor of *ajatella*, namely, human GROUPS, PARTICIPLES, and *että*-clauses as PATIENT. In contrast, for seven features the odds were against *ajatella*, and thus in favor of any one of the three other lexemes, these features being human GROUPS as AGENT, and NOTIONS, ATTRIBUTES, ACTIVITIES, forms of COMMUNICATION, INDIRECT QUESTIONS, and DIRECT QUOTES as PATIENT. In general, this comparison suggests that one cannot rely on a simple fit alone, as the different bootstrap sampling schemes reveal that potential variability, represented by the confidence intervals, is too broad for many of explanatory variables to be considered reliably and generally significant. Furthermore, it seems that the more rigorous the sampling scheme is, the more variability there is, thus reducing the number of effects assessed as significant, with the within-cluster sampling procedure producing the most stringent results. Finally, the addition of one explanatory variable to the model, representing an entirely different type of feature from the originally selected ones (extralinguistic vs.

morphological/syntactic/semantic), was observed to have a substantial impact on the weightings of the original variables. This underlines the importance of carefully considered variable selection, building upon a comprehensive understanding of the factors at work, which can be achieved by the combination of domain-specific knowledge of potential candidate types of variables and their selection through univariate and bivariate scrutiny.

Table 3.49. Comparison of the different fitting and sampling schemes of a binary logistic regression model of the selected same features (derived from Tables 3.45-3.48), with respect to the occurrence of *ajatella* against the rest. Significant positive association with *ajatella* is indicated by ‘+’, a negative positive association with ‘-’, and a nonsignificant result with ‘0’; results differing from the original single-round fit with the entire data are marked out in italics and with thicker border-lines.

Feature/Measure (<i>ajatella</i>)	Original model with single fit	Original model with simple bootstrap	Original model with within-cluster bootstrap	Original model + <i>medium</i> with single fit
Z_SG1	+	+	0	+
Z_SG2	0	0	0	0
Z_SG3	0	0	0	0
Z_PL1	+	+	0	+
Z_PL2	0	0	0	-
Z_PL3	+	+	0	+
SX_AGE.SEM_INDIVIDUAL	0	0	0	-
SX_AGE.SEM_GROUP	-	-	-	-
SX_PAT.SEM_INDIVIDUAL	+	0	0	+
SX_PAT.SEM_GROUP	+	+	+	+
SX_PAT.SEM_NOTION	-	-	-	-
SX_PAT.SEM_ATTRIBUTE	-	-	-	-
SX_PAT.SEM_STATE	-	0	0	-
SX_PAT.SEM_TIME	0	0	0	0
SX_PAT.SEM_ACTIVITY	-	-	-	-
SX_PAT.SEM_EVENT	0	0	0	0
SX_PAT.SEM_COMMUNICATION	-	-	-	-
SX_PAT.SEM_COGNITION	0	0	0	0
SX_PAT.SEM_LOCATION	0	0	0	0
SX_PAT.SEM_ARTIFACT	0	0	0	0
SX_PAT.INDIRECT_QUESTION	-	-	-	-
SX_PAT.DIRECT_QUOTE	-	-	-	-
SX_PAT.INFINITIVE	+	+	0	+
SX_PAT.PARTICIPLE	+	+	+	+
SX_LX_että_CS.SX_PAT	+	+	+	+

This process, which has been presented above for only one lexeme out of the four and with only a subset of explanatory variables to be included in the final model, with the different variants in sampling and validation, is exactly the same which will be applied to each of the component binary logistic models in the various heuristics for polytomous regression presented earlier in Section 3.4.3. In the full multivariate results to follow, with respect to the final set of explanatory variables selected as a

result of the univariate and bivariate analyses, I will only present the resultant odds ratios and the corresponding estimates of significance, starting with the simple fit of the model a single time with the entire data. I will then follow up with the assessment of the robustness of the effects by calculating confidence intervals, using both the simple bootstrap and the within-cluster scheme with writers/speakers as clusters. Finally, I will evaluate the effect of including the medium into the model.

3.4.6 Other possible or relevant multivariate methods

Other potential and relevant alternatives to logistic regression for multivariate analysis are the *probit* model, *discriminant analysis*, and *mixed-effects modeling*. In many respects, the probit model is similar to logistic regression, but the resultant parameters for a fitted probit model do not have a natural interpretation, thus rendering it less attractive (e.g., Fox 1997: 444-446; see also Agresti 2002: 246-247). Discriminant analysis is an older and once commonly used method especially in the case of polytomous outcomes, and it is simpler in terms of its calculation. However, it makes assumptions about the normality of the individual and joint distributions of the underlying variables which will not in practice hold, especially in the case of nominal variables. Even if these assumptions would be satisfied, regression analysis has been shown to be virtually as accurate as discriminant analysis, therefore indicating logistic regression as the more general analysis method (Harrell 2001: 217). Furthermore, discriminant analysis does not estimate instance probabilities *directly*, since in contrast to logistic regression it is based on the estimation of weights of predictor variables (the *X* in the regression formulas above) given some distribution of outcomes (the *Y* above).⁷⁹ What is more, these calculated parameter weights do not have a natural interpretation. In earlier similar linguistic studies, discriminant analysis has been used by Gries, for instance, for the analysis of the particle placement of phrasal verbs and the dative alternation in English, both having a dichotomous outcome (2003a, 2003b). Mixed-effects modeling⁸⁰ (Baayen et al., to appear 2008) represents rather a more advanced level of analysis than an equal alternative to logistic regression. In the case of this study, mixed-effects modeling would allow for incorporating, for instance, speaker/writer bias straightforwardly as a part of the actual statistical model, so that even speaker/writer-specific longitudinal effects are taken into account.

Furthermore, *classification and regression trees* (also known as *CART* models) and their extension *Random Forests*⁸¹ (Breiman 1995; Breiman and Cutler 2005) could also be an interesting supplement to compare the results of polytomous logistic regression with, as would *support vector machines* (SVM), various example-based rule-learning algorithms, and so forth. For instance, Gries (2003a) compared the prediction efficiency of discriminant analysis with a CART model. As was discussed above in Section 3.4.2, principal component analysis (PCA), as well as the older method of factor analysis (FA), or the latest modification, independent component analysis (ICA), could be used to cluster and reduce the overall number of variables.

⁷⁹ Estimates of outcome probabilities can, nevertheless, be derived from a discriminant model, but this requires inverting the model using Bayes' rule (Harrell 2001: 217).

⁸⁰ I am thankful to both of my external reviewers, Stefan Th. Gries and R. Harald Baayen, for reminding me of this method.

⁸¹ I am grateful to my external reviewer R. Harald Baayen for suggesting this method to me.

Moreover, this characteristic could also be used to study the overall relationships of the individual variables. Nevertheless, rather than comparing the many different ways of crunching the numbers, methods that instead aim at decreasing and compressing the complexity represented by multiple variables into a visual form would probably be the best complement to any numerical multivariate analysis, such as the polytomous logistic regression strategies presented above. Such methods include *correspondence analysis* (Lebart et al. 1998; see also, e.g., Agresti 2002: 382-384) and *self-organizing maps* (SOM), introduced by Kohonen (1995) as an offshoot of artificial neural networks, as well as *cluster analysis*.

With respect to Finnish, I have used correspondence analysis in my earlier work to study the distribution of morphological features among nouns (Arppe 2001) and verbs at various levels of granularity of semantic similarity (Arppe 2006b). This has included even the closest-knit synonym level, where my examples consisted of, among others, the four studied THINK lexemes (Arppe 2005a). Correspondence analysis is an attractive technique because it establishes for the items that it visually arranges a center with a surrounding periphery, reminiscent of the linguistic concept of prototypicality. Of the studied THINK lexemes, *ajatella* was placed closest to the visual origin, when the distributions of morphological features for it and the other three lexemes were taken into account. Another relevant example, employing self-organizing maps for the visualization of the collocational characteristics of a set of some of the most common Finnish verbs, has been undertaken by Lagus and Airola (2001). Though not concerning Finnish but highly relevant with respect to synonymy, cluster analysis has been applied to structure groups of near-synonymous Russian verbs denoting TRY and INTEND, building on their contextual (or *Behavioral*) profiles (Divjak and Gries 2006; Divjak 2006). One could easily extend these visual approaches to scrutinize relationships between the syntactic argument types, the semantic and structural of selected argument types, or any other sets of related variables (or all of them together, in accordance with Divjak and Gries), and the studied lexemes.

These visual methods mentioned above appear especially adept for determining the extent of semantic similarity between lexemes. Furthermore, the visual methods do build upon and thus contain precise numerical analysis, the results of which could be used to describe the associations of the lexemes and the features, as Divjak and Gries (2006) demonstrate. Nevertheless, such numerical data (e.g., *t-scores* and *z-scores* in the case of cluster analysis) lack the direct natural interpretation that logistic regression provides, in the form of odds for the explanatory variables and expected probabilities for the outcomes. Moreover, cluster analysis works on the aggregated proportions of the various features per outcome class, which, for example, constitute the (contextual) *behavioral profiles* in Divjak and Gries (2006: 36). Consequently, cluster analysis does not consequently take into consideration instance-wise co-occurrence patterns, that is, interactions among features (cf. Tables 2 and 4 in Gries and Divjak, forthcoming).⁸² Therefore, even though such visual analysis is strongly recommended in exploratory data analysis, as is indeed the case here (see, e.g., Hartwig and Dearing 1979), I have decided to exclude them from this study in order

⁸² Co-occurrences of selected features could, of course, be supplemented as separate, additional variables, though their selection should probably be prudent to do by hand, since the relative proportions of most co-occurrences would most probably be zero or close to it. Nevertheless, this would only cover the pairwise co-occurrences of features.

to be able to cover the selected numeric methods, presented in this Section 3, in a sufficiently thorough and comprehensive manner, and to retain some semblance of focus. However, in future studies, as so often is the case in science, it would seem most recommendable to combine both approaches and to capitalize on the advantages of each, by employing, firstly, visual methods for the determination of synonym groups, and, secondly, logistic regression analysis for describing the effects of the underlying contextual variables.

4 Univariate and bivariate results

4.1 Univariate analyses

4.1.1 General results

The array of various univariate statistical analyses, presented above in Section 3.2, for the distributions of singular features among the studied THINK lexemes were calculated using the R functions `explore.distributions(THINK.data, think.lex, "...")` and its subservient `singular.feature.distribution(THINK.data, think.lex, tag="...")`, and are presented in the data table `THINK.univariate` available in the `amph` data set. A selected subset of these results is presented in Appendix P. Furthermore, the intervening grouped-feature analyses have been calculated using the R function `multiple.feature.distribution(THINK.data, think.lex, tags=c(...))`.

In all, there were 477 contextual features or feature clusters in the final data which had at least the established overall minimum frequency (≥ 24) in the research corpus, and will thus constitute the major focus of scrutiny hereinafter. Of these, 378 (79.2%) exhibited a statistically significant ($P < 0.05$) overall heterogeneity in their distribution among the studied THINK lexemes, while their mean *Power* was 0.831 (s.d. 0.258) and the mean *Effect Size* $w = 0.0927$ (s.d. 0.0592). According to Cohen's (1992: 157, Table 1) proposals, such a mean *Effect Size* could barely be classified as *small* (for which the conventional minimum would be *Effect Size* = 0.10). For 123 (32.5%) of such features with overall significant distributions, the cellwise simplified abstracted results (+/-/0) as described in Section 3.2.1 above were exactly the same for all the four lexemes under consideration, regardless of whether the cell-wise (i.e., lexeme-wise) critical level was based on the minimum $X^2(df=1, \alpha=0.05)$ or equal to the one required for the overall (2x4) contingency table $X^2(df=3, \alpha=0.05)$, or scrutinized on the basis of the standardized Pearson residuals (with a critical level $|e_{\text{Pearson, standardized}}| \geq 2$). One should note, however, that such congruencies are a result of fortuitous combinations of an overall frequency and its distribution among the studied lexemes with respect to some particular features, rather than a systematic hierarchic relationship between these three criteria.

Out of all the theoretical 1512 (378·4) possibilities of feature-lexeme associations for features with overall significant distributions, there were 932 (61.6%) cellwise lexeme-specific significant associations (either '+' or '-') on the basis of the standardized Pearson residuals, 814 (53.8%) using the minimum $X^2(df=1, \alpha=0.05)$ value, and 542 (35.8%) with the conservative minimum $X^2(df=3, \alpha=0.05)$ value based on the overall table. Thus, the standardized Pearson residuals would appear to have overall the lowest threshold for suggesting a distinctive association, or disassociation, between some feature and an individual lexeme, with the minimum $X^2(df=1)$ trailing quite close behind. Consequently, I will stick to standardized Pearson residuals in the rest of the analyses to follow. This sensitivity is useful if one is after the smallest possible traces of distinctions among the studied lexemes; however, the down-side with such a low threshold is that it most probably is associated with a higher potential for refutation by other data or methods, in comparison to the more conservative measures.

With respect to the various measures of association, I had already selected Cramér's V as the only symmetric measure due to its simplicity and direct connection with the X^2 test of distributional homogeneity. Among the asymmetric measures, I have opted to retain, for the time, being the Goodman-Kruskal λ to provide continuity and a comparable reference point with Gries (2003a). However, the asymmetric Goodman-Kruskal τ as well as Theil's Uncertainty Coefficient U can both be considered to have better properties in comparison to the Goodman-Kruskal λ , as the former two take into account the entire distribution of possible outcome classes and not only the mode as is the case with the latter measure. Table 4.1 below presents the Pearson correlations of the values for all three of these association measures for the individual relationships of all 477 features exceeding the established frequency minimum threshold (≥ 24) with the studied THINK lexemes.

As can be seen among the asymmetric measures, the values of the two directions of the Goodman-Kruskal τ statistic, i.e., $\tau_{L[exeme]|F[eature]}$ and $\tau_{F[eature]|L[exeme]}$ correlate with each other to a very high extent, whereas there is almost no correlation between the two directions of the Goodman-Kruskal λ , i.e., $\lambda_{L[exeme]|F[eature]}$ and $\lambda_{F[eature]|L[exeme]}$, while the two directions of the Uncertainty Coefficient, i.e., $U_{L[exeme]|F[eature]}$ and $U_{F[eature]|L[exeme]}$, fall quite in the middle with moderate mutual correlation. Consequently, the feature-wise $\tau_{L|F}$ correlates strongly with both the corresponding $U_{L|F}$ and the lexeme-wise $U_{F|L}$, but the correlation between the corresponding lexeme-wise measures $\tau_{F|L}$ and $U_{F|L}$ is only moderate. Somewhat perplexingly, the symmetric Cramér's V correlates strongly with the feature-wise asymmetric measures $\tau_{L|F}$ and $U_{L|F}$ as well as the lexeme-wise $\tau_{F|L}$, but only moderately with the lexeme-wise $U_{F|L}$, and slightly less with either λ measures. This can be taken to underline the difference in the conceptual basis of Cramér's V against the asymmetric PRE measures, as it is a simple though effective normalization of the heterogeneity scrutinized with the X^2 statistic into the convenient [0,1] range, which is of practical use in ordering features but has little intrinsic meaning beyond that.

Table 4.1. The mutual Pearson correlations of the values of various nominal association measures based on the frequencies and distributions of all features exceeding the minimum overall frequency among the studied THINK lexemes.

```
cor (THINK.univariate[which (THINK.univariate [ ["freq" ] ] >=24) , c ("
cramers.v", "lambda.LF", "lambda.FL", "tau.LF", "tau.FL", "uc.LF", "
uc.FL" ) ], method="pearson")
```

Measures	Cramér's V	$\lambda_{L F}$	$\lambda_{F L}$	$\tau_{L F}$	$\tau_{F L}$	$U_{L F}$	$U_{F L}$
Cramér's V	1	0.561	0.454	0.936	0.937	0.944	0.629
$\lambda_{L F}$	-	1	0.0338	0.565	0.568	0.528	0.505
$\lambda_{F L}$	-	-	1	0.538	0.574	0.577	0.120
$\tau_{L F}$	-	-	-	1	0.970	0.985	0.581
$\tau_{F L}$	-	-	-	-	1	0.988	0.546
$U_{L F}$	-	-	-	-	-	1	0.568
$U_{F L}$	-	-	-	-	-	-	1

Furthermore, for the 477 individual features exceeding the minimum frequency threshold, the values of the lexeme-wise $\tau_{F|L}$ and $U_{F|L}$ measures are categorically greater than the corresponding feature-wise $\tau_{L|F}$ and $U_{L|F}$ values, and these differences are as a whole also statistically significant (one-sided paired t-test for $\tau_{F|L}$ and $\tau_{L|F}$: $t=-13.791$, $df=476$, $P<2.2e^{-16}$; for $U_{F|L}$ and $U_{L|F}$: $t=-22.30$, $df=476$, $P<2.2e^{-16}$);

however, in the case of the cruder measures $\lambda_{L|F}$ and $\lambda_{F|L}$ this asymmetry holds only in a minority of 14 (2.9%) cases and furthermore the differences are not overall significant (two-tailed paired t-test for $\lambda_{F|L}$ and $\lambda_{L|F}$: $t=0.542$, $df=476$, $P=0.706$). This is reflected also among the various statistics for the ranges of these measures presented in Table 4.2, as the lexeme-wise means of $\tau_{F|L}$ and $U_{F|L}$ are greater than those for the feature-wise $\tau_{L|F}$ and $U_{L|F}$, whereas the mean feature-wise $\lambda_{L|F}$ is slightly greater than the mean lexeme-wise $\lambda_{F|L}$, but in contrast to the other measures the standard deviation for $\lambda_{F|L}$ is of a magnitude greater than its mean.

As all three measures, λ , τ , and U , are conceptually based on the proportionate reduction of error (PRE), this result can be interpreted as indicating that a greater portion of the overall variation of the singular features among the studied lexemes is determined lexeme-wise rather than feature-wise in the singular-feature scrutiny, as is hypothesized on the basis of theoretical considerations in Appendix L. That is, knowing the lexeme increases the chances of guessing correctly the probability of whether a particular feature occurs with it or not more than predicting the probability of a lexeme knowing the feature. Therefore, these results would suggest that features are not fundamentally that monogamous with respect to which lexemes they occur with; rather, each of the lexemes may have its individual preference or dispreference with respect to a feature, and consequently more than one of the lexemes may have a similar preference. This may at least partially be attributed to the general setup of the singular-feature analysis, where there are more alternative categories available for lexemes than for features which are only considered in terms of their occurrence or nonoccurrence, the latter category which may bundle together a number of possible alternative (complementary) features logically related with the specific one under consideration.

Nevertheless, though the maxima for the lexeme-wise measures are moderately high at 0.19–0.21–0.27, the average values are considerably lower at 0.003–0.01–0.03, not to mention that the feature-wise means are of a magnitude lower at 0.00413–0.00413–0.00479, so overall neither the selected lexemes nor the contextual features by themselves can individually account for but a small portion of the observed usage. In the end, because the two Uncertainty Coefficient measures $U_{L|F}$ and $U_{F|L}$ exhibit real practical asymmetry in their value ranges, as is also evident in their density distribution for features exceeding the minimum frequency threshold in Figure 4.1, I will use them along with the symmetric Cramér's V in the later analyses below.

Table 4.2. The mean values, standard deviations, maxima and minima for Cramér's V , $U_{L|F}$, and $U_{F|L}$ for all features exceeding the minimum threshold frequency.

Association measure	Cramér's V	$\lambda_{L F}$	$\lambda_{F L}$	$\tau_{L F}$	$\tau_{F L}$	$U_{L F}$	$U_{F L}$
Mean	0.0927	0.00413	0.00346	0.00461	0.0121	0.00482	0.0373
Standard deviation	0.0592	0.00992	0.025	0.00688	0.0184	0.00712	0.0353
Maximum	0.433	0.107	0.274	0.0497	0.188	0.0566	0.208
Minimum	0.00718	0	0	0.00001	0.00005	0.00002	0.00050

Range and density of featurewise and lexemewise association values

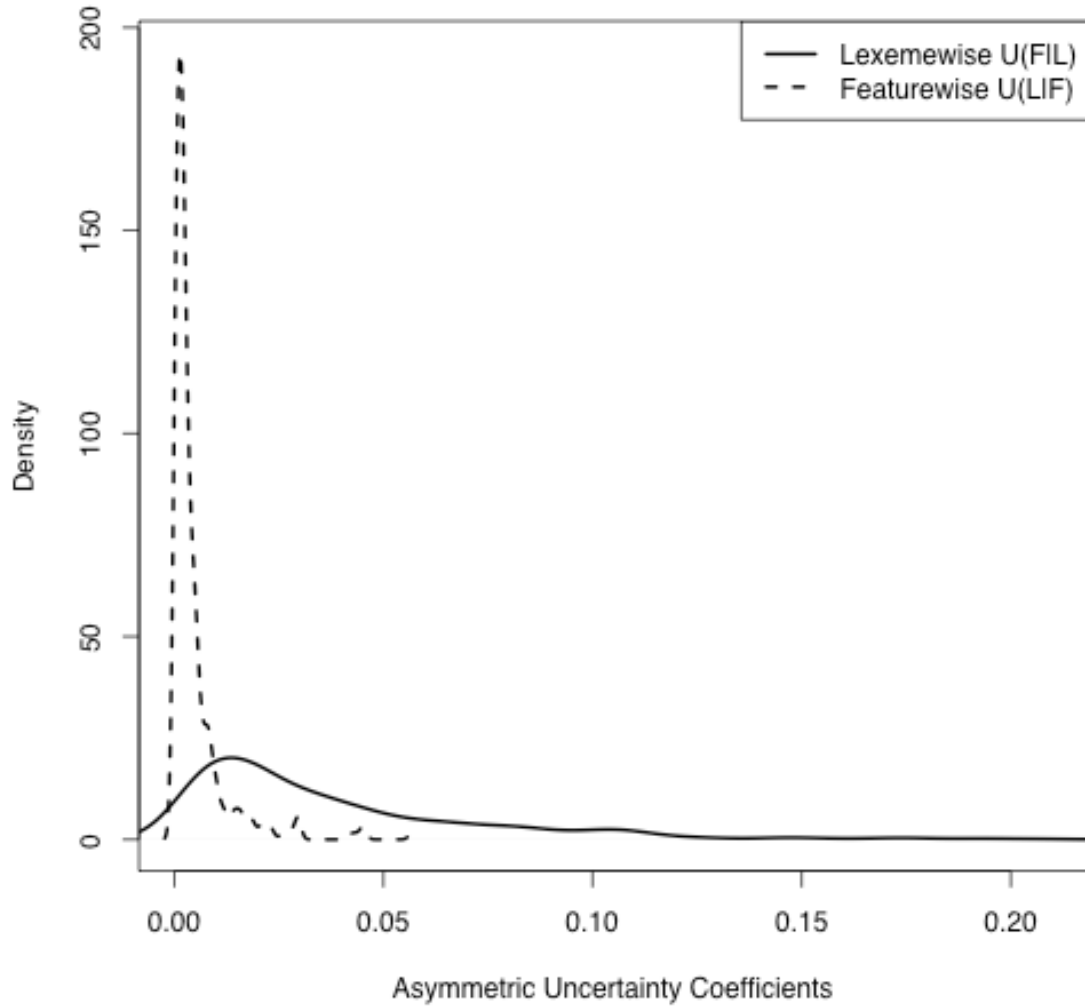


Figure 4.1. The distributions and densities of the two Uncertainty Coefficient measures $U_{L|F}$ and $U_{F|L}$ for features exceeding the minimum frequency threshold (≥ 24).

After scrutinizing the ranges and mean values of the various association measures we may first turn to which are the individual features that score the very highest with the selected measures. The 20 topmost features or feature combinations ranked by Cramér's V and both directions of the Uncertainty Coefficient $U_{L|F}$ and $U_{F|L}$ are presented in Table 4.3. In the first place we can notice that the topmost lexeme-wise ranked features have as a whole higher values (*minimum*=10.91) than the topmost feature-wise ranked features (*maximum*=0.0566), which is in line with the general results presented above. Secondly, there are clearly observable differences among the topmost sets of features, as only two are exhibited for all three measures, namely, an ACTIVITY and an INDIRECT QUESTION as a PATIENT. Furthermore, whereas there are as many as 19 shared features in the top 20 for both Cramér's V and the feature-wise $U_{L|F}$, the lexeme-wise $U_{F|L}$ has only one more feature in common with $U_{L|F}$, namely, quoted text (Z_QUOTE, particular to the newspaper subcorpus), in addition to topmost features common to all three measures.

Thirdly, looking at the individual features in Table 4.3, we can note that practically all types of features and feature combinations are present, be it morphological features of the node verb (i.e., INSTRUCTIVE case, Z_INS), morphological features of the entire verb-chain (e.g., FIRST and SECOND PERSON SINGULAR combined, Z_ANL_SG12), syntactic argument types alone (e.g., PATIENT, SX_PAT), lexemes as any syntactic arguments (e.g., the noun *työryhmä* ‘committee’, SX_LX_työ_ryhmä_N) or as specific syntactic arguments (the same noun as an AGENT, SX_LX_työ_ryhmä_N.SX_AGE), morphological features or parts-of-speech of syntactic arguments (e.g., NON-FINITE forms as PATIENT, SX_PAT.SX_NFIN, or a noun as PATIENT, SX_PAT.N), semantic and structural subtypes of syntactic arguments (e.g., the aforementioned ACTIVITY or INDIRECT QUESTION as PATIENT), as well as the extralinguistic categories of text type classifications (e.g., subforum on human relationships, Z_EXTRA_DE_ihmissuhteet), registers (e.g., newspaper text, Z_EXTRA_SRC_hs95), various relationships to attributive structures (e.g., usage in an attributive structure following a quoted passage, Z_POST_QUOTE), and even individual authors (i.e., author #948 in the SFNET newsgroup subcorpus, Z_EXTRA_AU_sfnet_948).

Table 4.3. The 20 topmost features per each of the three selected nominal associations, the symmetric Cramér’s V and the asymmetric Uncertainty Coefficients both lexeme-wise (U_{FIL}) and feature-wise (U_{LIF}), with respect to the distribution of the features among the studied THINK lexemes. Features in common for all three measure both in boldface and italics (and underlined); features common for Cramér’s V and U_{LIF} in boldface; features common for U_{LIF} and U_{FIL} in italics.

Cramér’s V (≥ 0.227)	U_{FIL} (≥ 0.109)	U_{LIF} (≥ 0.0208)
<u>SX_PAT.SEM_ACTIVITY</u> (0.433)	<i>SX_LX_työ_ryhmä_N.SX_AGE</i> (0.208)	<u>SX_PAT.SEM_ACTIVITY</u> (0.0566)
<u>SX_PAT.INDIRECT_Q...</u> (0.335)	<i>SX_LX_työ_ryhmä_N</i> (0.195)	<u>SX_PAT.INDIRECT_Q...</u> (0.0455)
<u>Z_EXTRA_SRC_hs95</u> (0.332)	<i>Z_POST_QUOTE</i> (0.190)	<u>Z_EXTRA_SRC_hs95</u> (0.0450)
<u>Z_EXTRA_SRC_sfnet</u> (0.332)	<u>SX_PAT.SEM_ACTIVITY</u> (0.176)	<u>Z_EXTRA_SRC_sfnet</u> (0.0450)
SX_PAT.N (0.332)	<i>SX_PAT.DIRECT_QUOTE</i> (0.174)	SX_PAT.N (0.0434)
Z_NON_QUOTE (0.330)	<i>SX_LOC.SEM_EVENT</i> (0.173)	Z_NON_QUOTE (0.0422)
SX_PAT (0.261)	<i>SX_MAN.SEM_GENERIC</i> (0.162)	SX_LX_että_CS.SX_PAT (0.0297)
SX_PAT.SX_SURF_NH (0.270)	<u>SX_PAT.INDIRECT_Q...</u> (0.152)	SX_LX_että_CS (0.0296)
SX_LX_että_CS.SX_PAT (0.267)	<i>SX_LX_niin_ADV.SX_MAN</i> (0.147)	SX_PAT.SX_SURF_CS (0.0296)
SX_LX_että_CS (0.267)	<i>SX_LX_näin_ADV</i> (0.146)	SX_PAT.CS (0.0296)
SX_PAT.SX_SURF_CS (0.266)	<i>SX_LX_näin_ADV.SX_MAN</i> (0.145)	SX_PAT.SX_SURF_NH (0.0294)
SX_PAT.CS (0.266)	<i>SX_LX_tapa_N.SX_MAN</i> (0.129)	SX_PAT (0.0281)
SX_AGE.N (0.248)	<i>SX_PAT.SX_NFIN</i> (0.127)	Z_ANL_SG12 (0.0274)
Z_EXTRA_DE_ihmis... (0.243)	<i>Z_INS</i> (0.122)	Z_EXTRA_DE_ihmis.... (0.0258)
SX_AGE.SEM_GROUP (0.241)	<i>SX_PAT.SX_PRON.SX_PHR_CLAUSE</i> (0.121)	Z_ANL_SGPL12 (0.0236)
Z_ANL_SG12 (0.240)	<i>SX_MAN.SEM_AGREEMENT</i> (0.115)	SX_AGE.N (0.0229)
Z_ANL_COVERT (0.232)	<i>SX_PAT.SX_FIN</i> (0.114)	Z_ANL_COVERT (0.0224)
SX_PAT.SX_SG (0.230)	<i>SX_PAT.SX_FIN.SX_PHR_CLAUSE</i> (0.113)	Z_POST_QUOTE (0.0221)
Z_ANL_SGPL12 (0.229)	<i>Z_EXTRA_AU_sfnet_948</i> (0.109)	SX_AGE.SEM_GROUP (0.0213)
SX_AGE.SX_SG (0.229)	<i>SX_PAT.PARTICIPLE</i> (0.109)	SX_PAT.SX_SG (0.0208)

However, it is probably of greater general interest to know how the various contextual feature categories as a whole fare on the average with respect to accounting for the occurrences of the scrutinized THINK lexemes. Such mean values according to the selected three association measures are presented in Table 4.4. In terms of Cramér’s V , the top-ranked category is the register/medium ($\kappa_{Cramer's V}=0.332$), followed by attributive structures (0.209), structural subtypes of syntactic arguments (0.135), verb-chain general morphological features (0.135), and syntactic argument types alone

(0.108). With respect to the asymmetric measures, taking the feature-wise perspective according to $U_{L|F}$, register/medium is again top-ranked (0.045), followed by attributive structures ($\chi^2_{U_{L|F}}=0.022$), structural subtypes of syntactic arguments (0.011), verb-chain general morphological features (0.009), and specific parts-of-speech as syntactic arguments (0.005). When considered lexeme-wise following $U_{F|L}$, the top two feature categories trade places, with attributive structures ranked the highest ($\chi^2_{U_{F|L}}=0.094$), followed by source/register (0.083), structural (0.073) as well as semantic (0.047) subtypes of syntactic arguments, and lexemes as specific syntactic arguments (0.047).

Interestingly, in comparison to the feature-wise ranking, verb-chain general morphological features appear lexeme-wise to have very little explanatory power ($U_{F|L}=0.009$ vs. $U_{L|F}=0.027$). Furthermore, although overall medium/register and association with attributive structures figure highest in terms of their association measure values, other extralinguistic feature categories such as repetition, author identity, and more fine-grained subsections within the two subcorpora have in contrast relatively very low mean values, as is the case also with morphological features in the non-node verb-chain context. The latter can be understood as a direct result from the fact that only a portion of the studied THINK lexemes occur in a verb-chain. If we look at the overall mean rankings of the various feature categories, register/medium, association with attributive structures, semantic and structural subtypes of syntactic arguments and verb-chain-general morphological features command the top ranks, and they are indeed the ones that will be included in the multivariate analysis later on.

Table 4.4. Mean values of selected association measures for various feature types; (number of features per type as well as ranking order of mean value per measure in parentheses).

Feature type	Cramér's V	$U_{L F}$	$U_{F L}$	Mean rank
Morphological feature of node verb (33)	0.106 (6)	0.00525 (9)	0.0312 (13)	9.3 (9)
Morphological feature in verb-chain (non-node) context (19)	0.0763 (16)	0.00260 (17)	0.0127 (18)	17.0 (18)
Morphological feature (anywhere) in verb-chain (25)	0.135 (4)	0.00893 (4)	0.0270 (15)	7.7 (6)
Syntactic argument type (alone) (19)	0.112 (5)	0.00596 (7)	0.0299 (14)	9.0 (8)
Lexeme as any syntactic argument (59)	0.0709 (17)	0.00269(16)	0.0405 (7)	13.3 (14)
Lexeme as specific syntactic argument (45)	0.0777 (15)	0.00318 (14)	0.0465 (6)	11.7 (13)
P-O-S feature of syntactic argument (39)	0.0979 (9)	0.00538 (8)	0.0318 (12)	9.3 (10)
Morphological feature of syntactic argument (132)	0.0866 (12)	0.00425 (12)	0.0350 (8)	10.7 (11)
Semantic subtype of syntactic argument (59)	0.0994 (8)	0.00520 (10)	0.0474 (5)	7.3 (5)
Structural subtype of syntactic argument (9)	0.135 (3)	0.0111 (3)	0.0734 (3)	3.0 (3)
Morphological or P-O-S feature of syntactic argument (141)	0.0892 (11)	0.00451 (11)	0.0343 (10)	10.7 (11)
Semantic or structural type of syntactic argument (68)	0.104 (7)	0.00598 (6)	0.0509 (4)	5.7 (4)
Author (15)	0.0634 (18)	0.00205 (18)	0.0336 (11)	15.7 (17)
Section (20)	0.0841 (13)	0.00409 (13)	0.0266 (16)	14.0 (15)
Source (2)	0.332 (1)	0.0450 (1)	0.0829 (2)	1.3 (1)
Repetition (7)	0.0796 (14)	0.00270 (15)	0.0181 (17)	15.3 (16)
Attributive structures (3)	0.209 (2)	0.0222 (2)	0.0936 (1)	1.7 (2)
All extra-linguistic features (47)	0.0954 (10)	0.00613 (5)	0.0343 (9)	8.0 (7)

As the number of variables that can be included in the multivariate analysis is in the order of tens rather than hundreds, some pruning of the possible features and feature combinations considered thus far will have to be undertaken at this stage to keep the following feature-specific analyses manageable. Overall, I will aim to select the most general of the available variables, matching the level of analysis presented earlier in Section 2.3.2 in conjunction with the scrutinies of the present lexicographical descriptions of the studied THINK lexemes. Thus, semantic as well as structural subtypes of syntactic arguments will be preferred over individual lexemes as specific arguments or as arguments in general, as the semantic classifications naturally cover a larger range of contexts, and the subtypes have been based on and are inspired by prominent individual lexemes. Indeed, 42 (71.2%) of the 59 individual lexemes as any syntactic argument exceeding the minimum frequency threshold were also matched in

the selected variable set by the same individual lexemes as specific syntactic arguments, further consisting of 34 (72.9%) cases in which the match was exclusively with one single type of syntactic argument. For the latter cases, the lexeme-specific (i.e., cellwise) preferences according to standardized Pearson residuals were exactly the same in 25 (73.5%) cases. And if we look for instance at the individual lexemes analyzed as the most frequent type of syntactic arguments, PATIENTS, their association with the semantic subtypes of this argument type is as high as $U_{Lexeme|Semantic_class}=0.941$.

Furthermore, as a sneak preview of the bivariate analyses to come, lexemes, as any type of syntactic arguments exceeding the minimum frequency threshold, exhibit a very high mean level of association with these same lexemes as particular syntactic argument types, with $U_{Lexeme|Lexeme+Syntactic_class}=0.901$. Nevertheless, subordinate clauses starting with the conjunction *että* ‘that’ as PATIENTS will be treated as a special exceptional case, since such structures are identified as a major distinct structural subtype in previous work (Pajunen 2001). In the case that the subclassification of some syntactic argument is skewed or scant, which applies for the less frequent syntactic argument types, the argument type alone will be used. Furthermore, while individual morphological features or parts-of-speech of various syntactic arguments are by themselves of interest, their great number renders their extended scrutiny impractical. Similarly, morphological features concerning the entire verb chain will be of greater use than the corresponding features specific to the node as well as the node-external components of the verb chain. Possible correlations among all the various features, however, will be of special interest in the later bivariate analyses, as this can uncover unexpected associations that may have been overlooked otherwise.

The comprehensive exposition and analysis category by category of individual singular-feature univariate results concerning the set of features selected above, as well as pertinent grouped-feature analyses, can be found in Appendix N. I will next move on to comparing these results with the existing lexicographical descriptions, as well as attempt to pull together the quite extensive assortment of features under a smaller set of *post hoc* generalizations.

4.1.2 Characterizations of the studied THINK lexemes on the basis of the univariate results

At first glance, the comprehensive run-through of the univariate results in Appendix N may look like a prolonged sequence of not very related details. This fault can be attributed to the exploratory approach chosen in the study, where there are no specific *a priori* hypotheses about the studied THINK lexemes, and where the central objective is to lay out and exemplify a methodological framework for studying the similarities and differences of lexemes within synonym sets. However, after wading through the extensive assortment of preference patterns, one can start to envision *post hoc* how they could be used to construct and support (or refute) more abstract characterizations concerning the core meanings of the selected THINK lexemes (Table 4.5).

At this point, I would venture firstly to designate *ajatella* as concerning temporally indefinite, continuous aspects of the cognitive process of “thinking”, undertaken by human beings individually, and often also concerning humans, or an intentional state.

In terms of agency, this characterization is matched by the preference exhibited by both human INDIVIDUALS and the FIRST PERSON for *ajatella*, and with respect to target/objective of the cognitive activity in both human INDIVIDUALS and GROUPS as PATIENT. Furthermore, holding a firm opinion or an enduring overall mental viewpoint is reflected in the GENERIC, FRAME, and AGREEMENT subtypes of MANNER associated with *ajatella*, and it is also most consistent with the neutral INDICATIVE mood indicating statements concerning states-of-affairs. Likewise, the fact that NEGATION is associated with *ajatella* also corresponds better with not generally having a stable opinion or viewpoint than the nonoccurrence of some individual fragment of thought in one's consciousness at a specific, delimited moment in time. Moreover, the propositions that the *että*-clause 'that' introduces are in a similar way temporally independent of the head verb itself in the main clause that such subordinate clauses are modifying.

In contrast, *mieltii* would be temporally more delimited and definite, though it would share the individualistic character of *ajatella*. Accordingly, in terms of individuality *mieltii* is preferred by the FIRST and SECOND PERSONS and SINGULAR number plus INDIVIDUAL AGENTS, coupled with ALONE as MANNER, and the IMPERATIVE mood, the latter which is typically associated with the SECOND PERSON addressing of other human individuals. The temporally restricted nature of *mieltii* is made evident in its preference by DIRECT QUOTES and INDIRECT QUESTIONS as well as nominal expressions of COMMUNICATION as PATIENTS, and explicit expressions of DURATION, whether LONG or SHORT, or the non-temporal but abstractly equivalent LITTLE subtype of QUANTITY, in addition to the potential for frequent repetition implied by OFTEN as FREQUENCY.

As for *pohtia*, it would be characterized by thinking undertaken by and as a group together, mostly concerning non-tangible, abstract notions. In this, the preference for *pohtia* by explicitly collective GROUP AGENTS is coupled neatly with positive associations with the human but impersonalized PASSIVE voice in addition to the more remote THIRD PERSON detached from the immediate discourse situation, as well as geographic LOCATIONS pertaining to human collective groups or collective activities such as EVENTS as LOCATIONS. With respect to the subtypes of PATIENTS, one can easily regard ATTRIBUTES as a subcategory or extension of abstract NOTIONS.

Finally, *harkita* would concern making decisions vis-à-vis actions, which would typically follow the actual cognitive process concerning taking such action.⁸³ Consequently, in addition to a preference for ACTIVITIES as PATIENT, *harkita* is preferred by AGAIN as FREQUENCY, implying the possibility of changing one's mind and future reconsideration, as well as by arguments denoting a REASON for contemplating the action in question, or a CONDITION necessary for making up one's mind. This conditionality of *harkita* is further reflected in its preference by both the CONDITIONAL mood and clause-adverbial META-arguments, the latter being typically used as hedges.

In the end, these lexeme-specific general characterizations can be seen to bear a resemblance to the Idealized Cognitive Models presented by Divjak and Gries (2006:

⁸³ This interpretation of *harkita* as possibly being overall future-orientated was first suggested to me by Professor Pentti Leino.

41-42), although they refer with the concept in question to small subclusters within a group of near-synonyms rather than to the individual lexemes themselves, to be precise. Furthermore, these general characterizations can be considered to represent the subconceptual, distinctive nuances among the studied THINK lexemes, at least with respect to their denotational dimension, within Edmonds and Hirst's (2002) three-level model of lexical choice, or equally well the different conceptualizations of similar events following Atkins and Levin (1995: 96), with the associated preferential contextual features as the explicit manifestations of these nuances or conceptualizations.

Table 4.5. Tentative hypotheses of the core semantic characteristics of the studied THINK lexemes and the associated contextual evidence (in the form of corpus-based relative [positive] preferences with respect to various features).

Lexeme	Semantic characterizations	Supporting contextual evidence (preferences)
ajatella	temporally continuous, individual in agency, and object, intentional state	Z_ANL_NEG, Z_ANL_IND Z_FIRST, SX_AGE.SEM_INDIVIDUAL SX_PAT.SEM_INDIVIDUAL SX_PAT.SEM_GROUP SX_LX_että_CS.SX_PAT SX_SOU SX_MAN.SEM_GENERIC SX_MAN.SEM_FRAME SX_MAN.SEM_AGREEMENT
miettiä	temporally definite, personal, individual in agency	Z_ANL_IMP Z_FIRST, Z_SECOND, Z_SING SX_AGE.SEM_INDIVIDUAL SX_PAT.SEM_COMMUNICATION SX_PAT.INDIRECT_Q..., SX_PAT.DIRECT_QUOTE SX_MAN.SEM_ALONE SX_QUA.SEM_LITTLE SX_DUR.SEM_LONG SX_DUR.SEM_SHORT SX_FRQ.SEM_OFTEN
pohtia	collective, impersonal in agency, non-concrete in object	Z_ANL_PASS, Z_ANL_THIRD, SX_AGE.SEM_GROUP SX_LOC.SEM_LOCATION SX_LOC.SEM_EVENT SX_PAT.SEM_NOTION SX_PAT.SEM_ATTRIBUTE
harkita	action as object, temporally in future	Z_ANL_KOND, SX_PAT.SEM_ACTIVITY SX_MAN.SEM_THOROUGH SX_FRQ.SEM_AGAIN SX_RSN SX_CND SX_META

With respect to *ajatella*, *pohtia*, and *harkita*, these characterizations and the associated preferences can also be seen to fit nicely with their more concrete, rural etymological origins. The INDIVIDUAL agency of *ajatella*, which might at first seem

somewhat odd considering the typically collective character that the etymologically underlying activities of *hunting* and *chasing* game nowadays have, receives a more fitting motivation when understood rather specifically in terms of *trapping*, a typically solitary kind of hunting even as it is still practiced today. One might also consider that the preference of *ajatella* in conjunction with human types of PATIENTS could be traced back to the general animacy of the objects of chasing/hunting, which specifically *pohtia* lacks in both its current usage and its origins. Though the preferred type of PATIENT for *pohtia* has changed from the concrete, that is, the grain and the chaff that are separated in *winnowing*, to the abstract, the collective nature of the original agricultural activity still clearly persists. Likewise, thoroughness and potential reconsideration as well as the (future) purpose-orientation now preferred by *harkita* are characteristics that one would still associate with the underlying activity of *trawling* with a dragnet. As a loan word adopted more or less in its current meaning, *miettiinä* alone among the selected THINK lexemes falls outside this historical continuum with respect to its contextual preferences.

4.1.3 Comparison of the univariate results with existing lexicographical descriptions

The univariate corpus-based results can be compared with the existing lexicographical descriptions in the two current dictionaries, namely, *Perussanakirja* (PS) and *Nykysuomen sanakirja* (NS), presented earlier in considerable detail in Section 2.3.2. In this, occurrence in the research corpus is only considered when the frequency of the feature in question has exceeded the minimum threshold value (≥ 24), and the granularity of the analysis is according to what was applied in the linguistic analysis of the example sentences provided in the dictionaries. Furthermore, syntactic argument types alone are considered only when no semantic and structural subtypes have been applied. Moreover, the extralinguistic features fall outside this comparison: firstly, because repetition is not applicable to the dictionary examples, and secondly the text types represented in the research corpus and those from which the dictionary example sentences are derived differ in their entirety, with the dictionaries consisting mainly of excerpts from well-known, established Finnish literature (as it was commonly conceived in the first half of the twentieth century).

Thus, the dictionaries and the research corpus contain a total of 102 distinct features, of which 42 (41.2%) are both mentioned in either dictionary and occur with sufficient frequency in the research corpus. These features could well be considered to have been thus demonstrated as characteristic and typical of the studied THINK lexemes on the whole and individually. However, the lack of overlap is considerable, as 39 (38.2%) features mentioned in either dictionary do not occur at all or with an insufficient (low) frequency in the research corpus, whereas 21 (20.6%) features exceeding the frequency threshold in the research corpus receive no mention among the examples in the dictionaries.

Details of the similarities and differences with respect to the feature set presented in the dictionaries and evident in the research corpus with sufficient frequency are given in Table 4.6 below. As can be seen, the only syntactic argument types that are evident in the research corpus but which are not represented by any semantic or structural subtype in the dictionary example sentences are GOAL and FREQUENCY. Thus, the

differences pertain rather to the granularity and selection of characteristic semantic and structural subtypes of these syntactic arguments among the studied THINK lexemes, as well as certain morphological features. It would also appear that the dictionaries have opted to include examples of rarer subtypes, at least in comparison to frequencies evident in the research corpus, such as BODY, ARTIFACT, and COMMUNICATION as subtypes of AGENT, and FAUNA, ARTIFACT, LOCATION, and COGNITION as subtypes of PATIENT.

Table 4.6. Details of the similarities as well as differences between the two dictionaries (PS and NS) and the corpus-based univariate singular-feature results; consideration of a feature as having sufficient occurrences in the research corpus requires meeting the minimum frequency threshold (≥ 24).

Feature category	Mention in either PS or NS as well as sufficient frequency in research corpus	No mention in either PS or NS but sufficient number of occurrences in research corpus (parenthesized features deliberately left as default values in the dictionary analyses)	Mention in PS or NS but no or insufficient occurrences in research corpus
MORPHO-LOGY	NEGATION, INDICATIVE, IMPERATIVE, PRESENT, PAST, FIRST, SECOND, THIRD, PLURAL, OVERT, COVERT, INFINITIVE1, INFINITIVE2, INFINITIVE3, PARTICIPLE1, PARTICIPLE2, TRANSLATIVE, INESSIVE, ILLATIVE, INSTRUCTIVE, CLAUSE EQ...	(AFFIRMATION), (ACTIVE), CONDITIONAL, (NOMINATIVE), GENITIVE, PARTITIVE, (SINGULAR/ PLURAL, CLITICS: <i>-kin/-pa</i>)	ESSIVE, ELATIVE, ABESSIVE
AGENT	INDIVIDUAL, GROUP	-	BODY, ARTIFACT, COMM...
PATIENT	INDIVIDUAL, NOTION, STATE, ATTRIBUTE, TIME, ACTIVITY, COMM..., INFINITIVE, INDIRECT_Q..., DIRECT_QUOTE, <i>että</i> -clause	GROUP, PARTICIPLE	FAUNA, ARTIFACT, LOCATION, COGNITION
SOURCE	NOTION	-	INDIVIDUAL
GOAL	-	-	INDIVIDUAL, NOTION, ATTRIBUTE, LOCATION
MANNER	GENERIC, FRAME, POSITIVE, THOROUGH, AGREEMENT (CONCUR), JOINT (ALONE)	NEGATIVE (\leftarrow SHALLOW)	CLARITY (\rightarrow POSITIVE), NOTION/ ATTRIBUTE, DIFFER, LIKENESS, ATTITUDE, SOUND
QUANTITY	MUCH, LITTLE	-	-
LOCATION	EVENT	LOCATION, GROUP	NOTION

TIME-POSITION	INDEFINITE	DEFINITE	-
DURATION	OPEN, SHORT, LONG	-	-
FREQUENCY	-	OFTEN, AGAIN	∅
VERB-CHAIN	NEGATED_AUX..., ADJACENT_AUX..., COMPLEMENT, PROPOSSIBILITY, IMPOSSIBILITY, PRONECESSITY, TEMPORAL, EXTERNAL, ACCIDENTAL	ABILITY (→ POSSIBILITY), NONNECESSITY, VOLITION (← TENTATIVE)	-
REASON	PURPOSE	-	REASON
CO-ORDINATED VERB	COGNITION	ACTION	THINK, MENTAL (→ PSYCHOL...)

The entire comparison of the lexeme-specific preference patterns for all the features considered here according to the corpus-based univariate singular-feature analyses, against mentions in example sentences in the existing two dictionaries, is presented in Table P.8 in Appendix P, while Table 4.7 below contains a summary of this comparison. We can again see that there is a substantial discrepancy. Of the 92 cases for which the results based on the research corpus have indicated a relative positive preference for one or more of the studied THINK lexemes with respect to a contextual feature, only 42 (45.7%) are mentioned in the examples in PS and 54 (58.7%) in NS, of which 40 (43.5%) are jointly apparent in both dictionaries, while as many as 37 (39.8%) remain unnoted in either dictionary. For the 222 instances of neutral feature-lexeme relationships according to the research corpus, 39 (17.6%) are noted in PS and 67 (30.2%) in NS, of which 21 (9.5%) are jointly mentioned, whereas 144 (64.9%) of such neutral cases are absent from the example sentences. With respect to the 100 dispreferences on the basis of the research corpus, 22 (22.0%) of such usages are nevertheless exemplified in PS and 37 (36.6%) in NS, of which 21 (21.0%) in both, while as many as 62 (62.0%) are not presented among the examples.

These results firstly reflect overall both the larger number of examples sentences provided in NS in comparison to PS, as well as the role of PS as a more concise successor to NS. Furthermore, though the dictionaries exemplify a greater part of the corpus-based lexeme-feature preferences, almost one-half of these remain unexemplified. Likewise, though the majority of lexeme-feature dispreferences do not occur among the example sentences, which is in accordance with the corpus-based results, a small but not altogether insignificant one-fifth of such dispreferred features are nonetheless provided as usage examples in the dictionaries.

Table 4.7. Comparison of lexeme-specific occurrences of features in the example sentences in the two dictionaries (PS and NS) against the preference patterns (+|0|-) derived with the univariate singular-feature analysis of the research corpus.

Lexeme-specific dictionary mention/ Preferences	Preference (+)	Neutrality (0)	Dispreference (-)	Σ
PS	42	39	22	103
NS	54	67	37	158
PS+NS	40	28	21	89
Ø	36	144	62	244
Σ	92	222	100	414

Table 4.8 below presents the specifics of features designated on the basis of the singular-feature analysis of the research corpus to have a positive preference for any of the studied THINK lexemes, but which nevertheless are not evident at all among the example sentences in either dictionary (PS or NS). For instance, we can see that GROUP as an AGENT in conjunction with *pohtia* has been omitted in both dictionaries, which does not correspond to the earlier results presented in Arppe and Järvikivi (2007b). Likewise, in the case of PATIENTS the positive preferences of the subtypes of human GROUPS, EVENTS, and PARTICIPLES with *ajatella* have been excluded in both dictionaries, as well as the preference of ATTRIBUTES, INDIRECT QUESTIONS, and DIRECT QUOTES in the same argument slot in conjunction with *pohtia*.

Table 4.8. Features designated with a positive preference for any of the THINK lexemes on the basis of the corpus-based singular-feature analysis which are not evident in the example sentences of either dictionary (PS and NS).

Contextual feature/Lexemes	<i>ajatella</i>	<i>mieltiä</i>	<i>pohtia</i>	<i>harkita</i>
MORPHOLOGY	-KIN	INFINITIVE4 -PA	PRESENT INFINITIVE3 INFINITIVE4 INESSIVE ILLATIVE	CONDITIONAL NOMINATIVE PARTITIVE
AGENT	-	-	GROUP	-
PATIENT	GROUP EVENT PARTICIPLE	-	ATTRIBUTE INDIRECT_Q DIRECT_Q	-
MANNER	GENERIC NEGATIVE	-	-	-
LOCATION	-	GROUP	LOCATION EVENT	-
TIME- POSITION	-	-	DEFINITE	-
DURATION	-	LONG SHORT	-	-
META (CLAUSE- ADVERBIAL)	-	-	-	UNSPECIFIED
VERB-CHAIN	-	NONNECESSITY VOLITION	TEMPORAL	-
CO- ORDINATED VERB	-	VERBAL	-	-

In contrast, Table 4.9 presents the specifics of features designated on the basis of the singular-feature analysis of the research corpus to have a dispreference for any of the studied THINK lexemes, but which nevertheless are presented among the example sentences in either one of the dictionaries or both (PS and/or NS). For instance, GROUPS as a subtype of AGENT is exemplified in the dictionaries in conjunction with *ajatella*, as are human INDIVIDUALS with *harkita*. Furthermore, with respect to PATIENTS, abstract NOTIONS, ATTRIBUTES, ACTIVITIES, and INDIRECT QUESTIONS are among the subtypes evident in the examples provided for *ajatella*, ACTIVITIES and (FIRST) INFINITIVES, not to mention *että*-clauses for *mieltiä*, and abstract NOTIONS and INDIRECT QUESTIONS for *harkita*. All of the aforementioned associations exemplified in the dictionaries are on the basis of the distributions of these features among the THINK lexemes in the research corpus analyzed as dispreferred usage (relative to the other selected THINK lexemes). In summary, the dictionaries appear to diverge substantially from what the corpus-based univariate results indicate as typical usage contexts of the studied THINK lexemes.

Table 4.9. Features designated with a dispreference for any of the THINK lexemes on the basis of the singular-feature analysis which are, however, nonetheless evident in the examples sentences of either dictionary (PS and NS).

Contextual feature/Lexemes	ajatella	miittää	pohtia	harkita
MORPHOLOGY	THIRD OVERT INFINITIVE3 INFINITIVE4 PARTICIPLE2 INESSIVE	PAST PASSIVE OVERT PARTICIPLE1 CLAUSE-EQ...	NEGATION	PRESENT PAST COVERT CLAUSE-EQ...
AGENT	GROUP	-	-	INDIVIDUAL
PATIENT	NOTION ATTRIBUTE ACTIVITY INDIRECT_Q	ACTIVITY INFINITIVE1 <i>että</i> -clause	-	NOTION INDIRECT_Q
MANNER	THOROUGH	-	-	GENERIC
TIME-POSITION	INDEFINITE	-	-	-
DURATION	OPEN LONG	-	-	-
VERB-CHAIN	ADJACENT- AUX... PRONECESSITY	-	NEGATED- AUX...	-
CO-ORDINATED VERB	VERBAL	-	-	-

4.2 Bivariate correlations and comparisons

4.2.1 Pairwise correlations of singular features

Comparing all the features in the original data table (as well as a few collapsed categories conceived of in the discussion of the univariate analyses) pairwise resulted in 227 475 feature pairings, of which 124 750 concern pairs with both features exceeding the minimum frequency threshold (≥ 24). In calculating the pairwise associations, the asymmetric Theil's Uncertainty Coefficient $U_{2|1}$ ⁸⁴ was selected, because the values for its two alternative directions remain in a two-feature setting (i.e., 2x2 contingency table) asymmetric in contrast to the Goodman-Kruskal τ , and thus allow us to evaluate whether either one of the features has a greater bearing on the other or vice versa. For the feature pairings satisfying the minimum frequency criterion, only 43 have a fully or close to perfect relationship, with $U_{2|1} \geq 0.99$, while an overwhelming majority of 100865 have practically no relationship at all with $U_{2|1} \leq 0.01$, as is also evident in the mean association value $U_{2|1} = 0.0163$ and the entire distribution of these values visualized in Figure 4.2, extremely skewed to zero on the left. Considering the intermediate range, 782 feature pairings have an unequivocally *strong* association with $U_{2|1} \geq 0.5$, while 2144 have at least a *moderate* relationship with $U_{2|1} \geq 0.2$, and 3905 at least a *weak* relationship with $U_{2|1} \geq 0.1$. For all of these pairings with at least a weak relationship the association is also statistically significant ($\forall U_{2|1} > 0.1 \Rightarrow P[U_{2|1}] < 0.05$). The full results of the pairwise comparisons have been calculated using the *R* function

```
singular.pairwise.association(cbind(THINK.data[THINK.univariate.tags.classified[,2]], THINK.data.extra),  
rbind(THINK.univariate.tags.classified,  
THINK.univariate.tags.extra.classified), compare="UC")
```

These are presented in the data table `THINK.bivariate` available in the `amph` data set, and those satisfying the minimum frequency requirement and with an asymmetric association value of at least $U_{2|1} > 0.1$ in the data table `THINK.bivariate.n_24.uc_.1` at the same location.

⁸⁴ In the notation used here and later on with respect to $U_{2|1}$ when used *alone*, *without* the corresponding $U_{1|2}$, this means that, in terms of the asymmetric Uncertainty Coefficient, any first mentioned feature (1) explains more of the variation of the associated second mentioned feature (2) than the other way around, i.e., $U_{2|1} \geq U_{1|2}$.

Range and density of pairwise association values

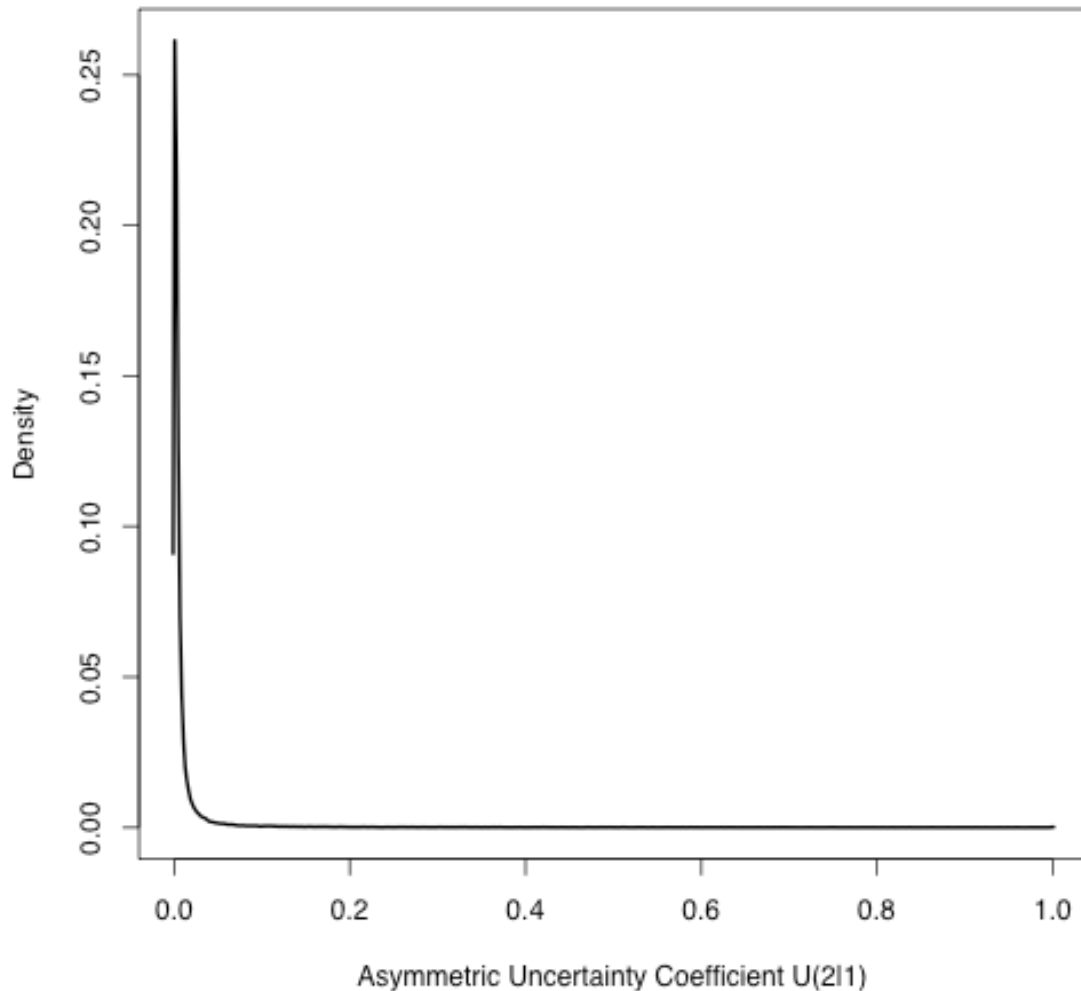


Figure 4.2. The distribution and density of the greater Uncertainty Coefficient values ($U_{2|1}$) for feature pairings both exceeding the minimum frequency threshold (≥ 24), with overall $n=123256$.

The topmost ten features in terms of their mutual association are presented in Table 4.10 below. As can be seen, the NON-FINITE and FINITE features have a perfect as well as complementary relationship, as their mutual $U_{2|1}$ and $U_{1|2}$ values are both exactly equal to one, but they have zero common occurrences. A similar perfect complementary relationship exists between the two features denoting the two subcorpora ($Z_EXTRA_SRC_hs95$ for the newspaper material and $Z_EXTRA_SRC_sfnet$ for the Internet newsgroup discussion data). We can also notice that other perfect associations follow from overlap in the underlying linguistic analysis scheme, for example, the relationship between being the first THINK lexeme in some text (Z_PREV_FIRST) and having no preceding THINK lexemes within the same text (Z_PREV_NONE), or from statistical (though not surprising) perfect overlap in that individual lexemes as syntactic arguments occur as only one particular argument type, for example, *edes* ‘even, at least’ as a clause-adverbial META-argument, or *joutua* ‘must’ as an ADJACENT AUXILIARY in the verb chain.

Table 4.10. The topmost ten feature pairs in terms of the Uncertainty Coefficient value, calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) so that $U_{2|1} \geq U_{1|2}$ always; $F_1 \equiv F_2 \sim F_1$ is logically equivalent to F_2 so that $F_1 \subset F_2$ and $F_1 \supset F_2$; $F_1 || F_2 \sim F_1$ is logically complementary throughout the entire data with F_2 so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$ and $\forall x(x \in F_1 \vee x \in F_2)$.

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
Z NFIN Z FIN	1	1	1973	1431	0
Z ANL THIRD≡Z ANL_SGPL3	1	1	1519	1519	1519
Z EXTRA_SRC sfnet Z EXTRA_SRC hs95	1	1	1654	1750	0
Z PREV_FIRST≡Z PREV NONE	1	1	2641	2641	2641
SX LX edes ADV.SX META⊂SX LX edes ADV	1	1	42	42	42
SX AGE.SX SG3≡SX LX hän PRON.SX AGE	1	1	91	91	91
SX LX joutua V.SX AAUX⊂SX LX joutua V	1	1	49	49	49
SX LX jälkeen PSP.SX TMP⊂SX LX jälkeen PSP	1	1	24	24	24
SX LX koskaan ADV.SX TMP⊂SX LX koskaan ADV	1	1	34	34	34
SX LX kuitenkin ADV.SX META⊂SX LX kuitenkin ADV	1	1	28	28	28

Instead of the above, I am rather interested in “surprising” associations that do not arise from the characteristics of the linguistic analysis scheme, that is, 1) co-occurrences of morphological or verb-chain specific features which are not logically mutually or unidirectionally implied in the underlying analysis, 2) correlations among different syntactic argument types and their respective semantic and structural subtypes, 3) correlations between node-specific or verb-chain general morphological features on the one hand and syntactic arguments and their semantic and structural subtypes on the other hand, and 4) correlations of extralinguistic features with each other or with node-specific or verb-chain general morphological features or syntactic arguments and their semantic and structural subtypes. For this purpose, the features were classified pairwise using the script `compare-and-classify-bivariate-tags`.

As was noted in the presentation of bivariate methods above in Section 3.2.1, in addition to the explicit meaning of PRE measures, such as the Uncertainty Coefficient $U_{2|1}$, as the variation of one feature which is explained, or in more modest terms, accounted for without necessarily assuming direct causality, by the occurrences of the other feature, any verbal interpretations of such measures on the basis of some threshold values are always arbitrary. Nevertheless, I will apply such generally suggested gradings so that when $U_{2|1} > 0.5$, I will consider the association among the pair unequivocally *strong* enough to allow for only one of the features to be included in the following multivariate analysis. The actual selection or rejection of variables, however, will be undertaken later in Section 5.1 as a prelude to the multivariate analysis.

As the overall number of contextual features which are considered in this study to have a bearing on the usage of the studied THINK lexemes is relatively large, when an overwhelming majority of the pairwise associations are practically null, even the smaller but nonzero values graded as *moderate*, i.e., $U_{2|1} > 0.2$, may, due to their relative infrequency, turn out to be of greater interest than in some other circumstances. As in the case of pairings of syntactic arguments and their subtypes, there are very few moderate or strong relationships, I will even scrutinize *weak* associations with $U_{2|1} > 0.1$. Indeed, if we recall the univariate results with respect to

the associations of individual features with the studied THINK lexemes (in Table 4.4 in Section 4.1.1), we may note that the mean values as well as the distribution ranges were overall relatively low overall, at $\bar{x}(U_{F|L})=0.0373$, $\min(U_{F|L})=0.00050$, and $\max(U_{F|L})=0.208$ lexeme-wise, and even less feature-wise, at $\bar{x}(U_{L|F})=0.00482$, $\min(U_{L|F})=0.00002$, and $\max(U_{L|F})=0.0566$. The complete results of the bivariate comparisons for the selected feature pairings meeting the above criteria are presented in Appendix Q, of which the relevant results are presented here below.

We can start off with the associations of the node-specific morphological features, presented in Table 4.10. Several of these associations remind us of the logically complementary binary relationships in the underlying analysis scheme, namely, NON-FINITE (Z_NFIN) vs. FINITE (Z_FIN) forms, and ACTIVE (Z_ACT) vs. PASSIVE (Z_PASS) voice. Perhaps the most relevant relationships here are the very strong ones of the SECOND INFINITIVE with the INSTRUCTIVE case and the THIRD INFINITIVE with the ILLATIVE case, with $U_{2|I}=0.866$ and $U_{3|I}=0.748$, respectively, for which both the specific type of infinitive explains slightly more of the occurrence of the particular morphological case than vice versa. In my judgement, the former association arises from an idiosyncratic form based nearly always on *ajatella*, namely, *ajatellen*, roughly corresponding to ‘thinking about [something], with [something] in mind’. The latter is associated with the obligatory government required by some types of auxiliary verbs, for example, *sai/ryhtyi ajattelemaan* ‘got [someone] to think/[someone] started to think’. These relationships among particular types of infinitives and cases will motivate a grouped-feature analysis later on in Section 4.2.2 concerning the general types of NON-FINITE forms, that is, the five infinitives and the two participles, on the one hand and the nominal cases on the other.

Furthermore, among the strongest relationships, we can notice that SINGULAR *number* (associated by definition with NON-FINITE forms) is also clearly linked with the THIRD INFINITIVE, which is also coupled by the lesser but still relatively high association of this feature with the ILLATIVE case, suggesting an overall close relationship for this particular feature trio. However, despite such specific strong associations neither number feature appears to have a more significant general role among all the node-specific features considered here.

Continuing downwards, we may further see that the INDICATIVE mood is somewhat more associated with the PRESENT rather than the PAST tense. With respect to person-number features, the INDICATIVE mood has a moderate association with the THIRD PERSON SINGULAR as well as the FIRST PERSON SINGULAR and the THIRD PERSON PLURAL, whereas the IMPERATIVE mood is associated with the SECOND PERSONS SINGULAR and PLURAL. The latter relationship is quite expected on the basis of the prototypical use of the IMPERATIVE mood to convey commands, exhortations, and requests to other persons in the immediate (discourse) context, which also explains the association of the clitic *-pa* ‘but, now’ with this mood, as it can be used to hedge and soften an expression, for example, *mietipä* ‘but (now) think’. In contrast to the two afore-mentioned moods, the CONDITIONAL does not appear to have any substantial associations. Finally, we may make note of the last association barely exceeding the preset minimum threshold $U_{2|I} \geq 0.2$, namely, SECOND PARTICIPLE with the TRANSLATIVE case, which in my judgements stems from the ACCIDENTAL verb-chain construction also involving some form of *tulla* ‘come’ as an auxiliary verb, for example, *tulin ajatelleeksi* ‘I came/happened to think of’.

Table 4.11. The pairwise associations of the node-specific morphological features considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{2|1} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) and presented always so that $U_{2|1} \geq U_{1|2}$; $F_1 > F_2 \sim U_{2|1} > U_{1|2}$; $F_1 \subset F_2 \sim F_1$ is a logical subset of F_2 ; $F_1 || F_2 \sim$ logical complementarity throughout the entire data so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$ and $\forall x(x \in F_1 \vee x \in F_2)$; $F_1 | F_2 \sim F_1$ is logically pairwise disjoint with F_2 so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$; $F_1 \neq F_2 \sim F_1$ is logically multiply disjoint within a set of related features $\cup(F_1, \dots, F_n)$ so that $F_1 \not\subset \cup(F_2, \dots, F_n)$ and $\cup(F_2, \dots, F_n) \not\subset F_1$; associations covered more generally by some other(s) or otherwise considered less informative in (parentheses).

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
Z_NFIN Z_FIN	1	1	1973	1431	0
Z_INF2>Z_INS	0.866	0.75	166	137	137
Z_INF3>Z_ILL	0.748	0.676	309	267	253
(Z_SG>Z_INF3)	0.504	0.297	720	309	306
Z_IND>Z_PRES	0.474	0.423	1272	943	883
Z_SG2>Z_IMP	0.454	0.403	171	146	106
Z_SG>Z_ILL	0.418	0.223	720	267	255
(Z_IMP>Z_PA)	0.363	0.18	146	59	43
Z_IND>Z_PAST	0.353	0.19	1272	389	389
Z_IMP>Z_PL2	0.347	0.153	146	51	37
Z_IND>Z_SG3	0.328	0.209	1272	509	488
(Z_SG>Z_TRA)	0.285	0.045	720	54	53
Z_ACT Z_PASS	0.272	0.176	1624	561	0
(Z_ACT \subset Z_FIN)	0.256	0.252	1624	1431	1163
(Z_ACT \neq Z_NFIN)	0.256	0.252	1624	1973	461
Z_IND>Z_SG1	0.222	0.087	1272	248	234
Z_PCP2>Z_TRA	0.214	0.044	454	54	42
Z_IND>Z_PL3	0.2	0.058	1272	164	156

Moving on to the morphological features concerning the entire verb chains of which the studied THINK lexemes form part, we can again see in Table 4.12 among the strongest relationships the association of SECOND person forms with the IMPERATIVE mode, in concordance with the node-specific results above; this relationship is mirrored in the SECOND PERSONS SINGULAR and PLURAL features individually, though with a lesser association value. In contrast, the INDICATIVE mood does not exhibit any of the person-number associations related to the node-specific ones. Among themselves, however, the three frequent moods have relatively substantial associations, a practical example of multicollinearity arising from partial redundancy among mutually exclusive categories (remember, e.g., Cohen et al. 2003: 311). This factor is also present in the disjoint associations of AFFIRMATIVE polarity and INDICATIVE mood as the most frequent of their categories among finite verb-chains and CLAUSE-EQUIVALENT forms which can neither be marked for polarity nor mood, the SINGULAR number and THIRD person with the PASSIVE voice, as well as the corresponding positive association of the SINGULAR number with the ACTIVE voice, which must be taken into account in the final selection of variables in the multivariate analysis.

It is noteworthy that the PLURAL number does not exhibit any substantial associations with any other features other than complementarity with the SINGULAR number, a relationship which expectedly also applies for AFFIRMATIVE and NEGATIVE polarity,

ACTIVE and PASSIVE voice, as well as OVERT vs. COVERT manifestations of AGENTS/subjects. With respect to the latter feature pair, we may also note that a COVERT AGENT has a moderate association with the IMPERATIVE mood, while an OVERT AGENT has a similar relationship with ACTIVE voice in general and THIRD PERSON PLURAL in particular, which all arise from the conventions of standard written Finnish: as in English an AGENT/subject may be omitted in commands and requests expressed with the IMPERATIVE mood, while among the different types of AGENTS the THIRD person forms, especially in the PLURAL, are more of than not explicitly expressed in Finnish. Lastly, SINGULAR (nominal) number has a weak association with CLAUSE-EQUIVALENT forms.

Table 4.12. The pairwise associations of the verb-chain general morphological features considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{2|1} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) and presented always so that $U_{2|1} \geq U_{1|2}$; $F_1 > F_2 \sim U_{2|1} > U_{1|2}$; $F_1 \subset F_2 \sim F_1$ is a logical subset of F_2 ; $F_1|F_2 \sim F_1$ is logically pairwise disjoint with F_2 so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$; $F_1 \neq F_2 \sim F_1$ is logically multiply disjoint within a set of related features $\cup(F_1, \dots, F_n)$ so that $F_1 \not\subset \cup(F_2, \dots, F_n)$ and $\cup(F_2, \dots, F_n) \not\subset F_1$; associations covered more generally by some other(s) or otherwise considered less informative in (parentheses); ‘?’ indicates (minor) inconsistency in underlying analysis scheme.

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
Z ANL AFF≠Z PHR CLAUSE	0.623	0.48	2573	521	0
Z ANL SECOND>Z ANL IMP	0.585	0.342	320	152	147
Z ANL IND≠Z PHR CLAUSE	0.516	0.362	2386	521	0
Z ANL SING<?Z ANL ACT	0.509	0.47	1962	2306	1918
Z ANL AFF Z ANL NEG	0.471	0.259	2573	310	0
Z ANL ACT Z ANL PASS	0.445	0.279	2306	457	0
Z ANL IND≠Z ANL KOND	0.378	0.174	2386	275	0
(Z ANL SG2>Z ANL IMP)	0.366	0.25	256	152	111
Z ANL OVERT Z ANL COVERT	0.36	0.352	1314	1218	0
(Z ANL SGPL12>Z ANL IMP)	0.343	0.113	829	152	150
(Z ANL SING≠Z ANL PASS)	0.329	0.191	1962	457	0
Z ANL IND≠Z ANL IMP	0.309	0.092	2386	152	0
Z ANL SING Z ANL PLUR	0.304	0.158	1962	386	0
Z ANL COVERT>Z ANL ACT	0.288	0.278	1218	2306	1217
(Z ANL COVERT>Z ANL SGPL12)	0.273	0.232	1218	829	682
(Z ANL COVERT>Z ANL SG12)	0.258	0.202	1218	705	592
(Z ANL IMP>Z ANL PL2)	0.231	0.118	152	64	36
Z ANL SING>Z PHR CLAUSE	0.227	0.143	1962	521	44
Z ANL THIRD≠Z ANL PASS	0.222	0.128	1519	457	0
(Z ANL SING>Z ANL COVERT)	0.221	0.211	1962	1218	1109
Z ANL COVERT>Z ANL IMP	0.221	0.062	1218	152	147
Z ANL OVERT>Z ANL PL3	0.216	0.088	1314	262	247

Turning to the relationships among different syntactic arguments and their subtypes presented in Table 4.13, we may firstly note that their overall degree of mutual association is relatively lower than that which was evident for the node-specific and verb-chain general morphological features above; consequently, I have lowered the minimum threshold for inclusion in the considerations here to $U_{2|1} > 0.1$. In general, I interpret this to follow from the extensive pairwise combinatorial possibilities arising from the 19 syntactic argument types and their close to 70 semantic and structural

subtypes, which have exceeded the minimum frequency thresholds, not to mention the less frequent subtypes also evident in the research corpus. The only argument pair which exhibits a relatively strong relationship concerns CO-ORDINATED VERBS and CO-ORDINATED CONJUNCTIONS, which is exactly what one could expect; that this association is not entirely perfect arises from the few instances where a lexical coordinated conjunction has been omitted, in which case punctuation, typically a comma, is used instead.

Among the weaker relationships we may note that several combinations of PATIENT and MANNER arguments rise above the rest, specifically PATIENT arguments in general with the GENERIC, AGREEMENT, and CONCUR subtypes of MANNER. Consequently, I will conduct a pairwise grouped-feature analysis of the subtypes of both PATIENT and MANNER arguments in Section 4.2.2 below. Furthermore, we can see with respect to verb chains that their less frequent components, non-adjacent non-negation auxiliaries (SX_CAUX) as well as nominal complements (SX_COMP) are somewhat moderately associated with occurrences of adjacent auxiliaries (SX_AAUX). Moreover, nominals in general denoting TIME as TIME-POSITIONS as well as overall DEFINITE expressions of this same syntactic argument exhibit a weak association with an EVENT as a LOCATION. With a similar low degree of association, AGENT arguments in general are linked with GROUPS as LOCATION, and arguments denoting REASON with verb-chains expressing an EXTERNAL cause. These last mentioned relationships are all intuitively plausible ones, that is, it would seem natural to express an explicit point in time in conjunction with an event, as well as to have a reason causally followed by some consequent action.

Table 4.13. The pairwise associations of the different syntactic arguments and their semantic and structural subtypes considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{21} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 (U_{21}) and vice versa (U_{12}) and presented so that always $U_{21} \geq U_{12}$; $F_1 > F_2 \sim U_{21} > U_{12}$; $F_1 < F_2 \sim F_1$ is a logical subset of F_2 ; ‘?’ indicates (minor) inconsistency in underlying analysis scheme. Furthermore, correspondences arising among the various overlapping analysis schemes concerning verb-chains (SX_VCH.SEM_XXX, SX_XAUX, Z_ANL_XXX) or nonoccurrences of particular syntactic arguments (SX_XXX.SEM_NIL) mirroring a correspondence with a positive occurrence of the same argument (SX_XXX) are mostly ignored as noninformative.

Feature pair	U_{21}	U_{12}	n_1	n_2	n_{common}
SX CV<?SX CC	0.837	0.761	190	167	163
SX PAT>SX MAN.SEM GENERIC	0.262	0.082	2812	113	17
SX PAT>SX MAN.SEM AGREEMENT	0.215	0.034	2812	48	7
SX AAUX>SX CAUX	0.213	0.054	1271	134	131
SX PAT>SX MAN.SEM CONCUR	0.186	0.018	2812	26	4
SX AAUX>SX COMP	0.162	0.049	1271	171	154
SX TMP.SEM TIME>SX LOC.SEM EVENT	0.138	0.053	119	36	15
SX TMP.SEM DEF...>SX LOC.SEM EVENT	0.131	0.041	158	36	16
SX AGE>SX LOC.SEM GROUP	0.124	0.018	2537	56	12
SX VCH.SEM EXTERNAL>SX RSN	0.113	0.1	79	68	20
SX MAN>SX PAT	0.121	0.118	616	2812	326

After scrutinizing several individual feature categories separately we may move on to observing whether these feature sets exhibit any inter-category relationships. When pitting node-specific morphological features against syntactic argument types and

their semantic and structural subtypes, presented in Table 4.14, a majority of the at least moderate associations concern the structural components or semantic classifications of verb-chains. Firstly, a notable and very strong relationship exists between the TRANSLATIVE case and ACCIDENTAL verb-chains, in which the former feature is one of the three explicit component exponents of the latter, more semantic characterization for their co-occurrence, in which the other two members are *tulla* ‘come’ as an (adjacent) auxiliary and the PAST PARTICIPLE form of the node THINK verb. Indeed, this latter feature has a lesser, moderate association with the verb-chain subtype in question.

A similar relationship applies for the ILLATIVE case and the THIRD INFINITIVE with verb-chains expressing an EXTERNAL cause. Furthermore, verb-chains consisting of a THINK lexeme in the FIRST INFINITIVE are moderately associated with modality concerning POSSIBILITY, whether (positive) PROPOSSIBILITY or (negative) IMPOSSIBILITY, or to a lesser extent non-positive NECESSITY in general, i.e., SINENECESSITY. Moreover, immediately adjacent auxiliary verbs are associated foremost with FIRST INFINITIVE forms, the other possible alternative being the THIRD INFINITIVE, for which the strongest association has been mentioned above. The FIRST INFINITIVE is also linked with verb-chain nominal COMPLEMENTS in general and particularly those denoting abstract NOTION. In contrast to the aforementioned two types of infinitives, the SECOND INFINITIVE, shown to be connected above in a close relationship with the INSTRUCTIVE case which is evident also here, is associated with actual syntactic arguments, having a negative though not perfect association with AGENTS, but a positive one with EVENTS as PATIENTS.

Last among the positive associations we may note THIRD PERSON SINGULAR in conjunction with DIRECT QUOTES as PATIENTS, and a near equivalence between a negative auxiliary (SX_NAUX) and a negated form of the node verb (Z_NEG), the latter which is dictated by the principles of proper Finnish grammar. Finally, concerning the type of “inverse” multicollinearity noted earlier above, we can see that *not* having any indication of modality in the verb chain has a moderate negative association with the FIRST INFINITIVE, as is the case also between syntactic AGENTS in general, and to a lesser extent human INDIVIDUALS in particular, and the PASSIVE voice. That the occurrence of a syntactic AGENT nevertheless has some co-occurrence with the PASSIVE voice is due to certain CLAUSE-EQUIVALENT or neccessive constructions,

e.g., *asiaa*_{PATIENT+NOTION}

*harkittuani*_{PARTICIPLE2+PASSIVE+(PARTITIVE)+SG1,COVERT} ‘having considered the matter’ and *minun*_{AGENT+INDIVIDUAL+FIRST+GENITIVE} *on*_{A-AUX} *harkittava*_{PARTICIPLE1+PASSIVE...} ‘I must consider ...’, which morphologically employ the intuitively contradictory alternative (i.e., PASSIVE) among the two possible alternatives of voice.

Table 4.14. The pairwise associations of node-specific morphological features on the one hand and syntactic arguments and their semantic and structural subtypes on the other hand which have been considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{2|1} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) and presented so that always $U_{2|1} \geq U_{1|2}$; $F_1 > F_2 \sim U_{2|1} > U_{1|2}$; $F_1 \subset F_2 \sim F_1$ is a logical subset of F_2 ; $F_1 \equiv F_2 \sim F_1$ is logically equivalent to F_2 so that $F_1 \subset F_2$ and $F_1 \supset F_2$; $F_1 \neq F_2 \sim F_1$ is logically multiply disjoint within a set of related features $\cup(F_1, \dots, F_n)$ so that $F_1 \not\subset \cup(F_2, \dots, F_n)$ and $(F_2, \dots, F_n) \not\subset F_1$; associations covered more generally by some other(s) or otherwise considered less informative in (parentheses); ‘?’ indicates (minor) inconsistency in underlying analysis scheme. Furthermore, correspondences arising among the various overlapping analysis schemes concerning verb-chains (SX_VCH.SEM_XXX, SX_XAUX, Z_ANL_XXX) or nonoccurrences of particular syntactic arguments (SX_XXX.SEM_NIL) mirroring a correspondence with a positive occurrence of the same argument (SX_XXX) are mostly ignored as noninformative.

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
Z TRA \supset SX VCH.SEM ACCIDENTAL	0.807	0.684	54	44	42
(SX NAUX \equiv ?Z NEG)	0.514	0.24	314	111	106
Z ILL $>$ SX VCH.SEM EXTERNAL	0.439	0.176	267	79	71
Z INF3 $>$ SX VCH.SEM EXTERNAL	0.421	0.153	309	79	72
SX VCH.SEM NILMODALITY $>$ Z INF1	0.394	0.359	2572	695	132
SX AGE \neq Z PASS	0.393	0.31	2537	561	73
(SX AGE.SEM NIL \equiv Z PASS)	0.393	0.31	867	561	488
Z SG3 $>$ SX PAT.DIRECT_QUOTE	0.386	0.14	509	120	113
SX AAUX $>$ Z INF1	0.36	0.276	1271	695	649
Z PCP2 \supset SX VCH.SEM ACCIDENTAL	0.357	0.063	454	44	43
Z INF1 $>$ SX VCH.SEM POSSIBILITY	0.317	0.206	695	347	285
Z INF1 $>$ SX VCH.SEM NILPOSSIBILITY	0.317	0.206	695	3057	410
(SX AGE.SEM INDIVIDUAL \neq Z PASS)	0.296	0.207	2251	561	62
Z INF1 $>$ SX COMP.SEM NOTION	0.281	0.048	695	58	56
SX AGE $>$ Z INS	0.271	0.081	2537	137	8
(Z INF1 $>$ SX VCH.SEM PROPOSSIBILITY)	0.263	0.142	695	264	212
Z INF1 $>$ SX COMP	0.244	0.096	695	171	142
(Z INF1 $>$ SX VCH.SEM IMPOSSIBILITY)	0.232	0.053	695	83	73
SX AGE $>$ Z INF2	0.227	0.078	2537	166	21
(SX AGE.SEM INDIVIDUAL $>$ Z INS)	0.224	0.059	2251	137	5
Z INS $>$ SX PAT.SEM EVENT	0.218	0.063	137	29	17
Z INF1 $>$ SX VCH.SEM SINENECESSITY	0.211	0.036	695	57	50

Turning to the interrelationship between verb-chain general features and syntactic arguments and their semantic subtypes, presented in Table 4.15, the strongest association indicates the practically equivalent relation between an explicit negative auxiliary verb (SX_NAUX) and overall negative polarity, i.e., NEGATION (Z_ANL_NEG), of the verb chain. To a lesser extent, NEGATION is also moderately associated with the more general SINENECESSITY and its more particular subtype NONNECESSITY, and IMPOSSIBILITY, that is, in essence all nonpositive subtypes of modality applicable for verb-chains, though the frequencies of common occurrences show that this aspect can also sometimes be conveyed *without* the explicit prototypical NEGATION feature, with, for example, a nominal complement denoting impossibility or the like, i.e., *on*_{ADJACENT_AUXILIARY} *mahdotonta*_{COMPLEMENT} *ajatella* ‘[it] is impossible to think’.

Other notable relationships are the positive strong association of ACTIVE voice, and the moderate ones of both SINGULAR number and THIRD person, all with a syntactic AGENT as an argument. The inverse counterpart of these is the quite strong negative association of syntactic AGENTS in general, and to somewhat lesser extent its INDIVIDUAL subtype in specific, with the PASSIVE voice. Furthermore, tendencies concerning the explicitness of the subject/AGENT become now evident, as COVERTness has a moderate association with INDIVIDUAL AGENTS and to a lesser extent with AGENTS in general, whereas OVERTness has a moderate association with GROUPS as AGENT and DIRECT QUOTES as PATIENT, that is, in conjunction with attributive constructions. This is concordant with the fact that the FIRST and SECOND PERSON both SINGULAR and PLURAL pronouns which may be omitted are classified as INDIVIDUAL AGENTS in this study, manifested in the association between these features just exceeding the minimum threshold. Finally, PASSIVE voice in the verb chain is moderately associated with GROUPS as LOCATION arguments, in which case the LOCATION argument may in fact be considered to represent some characteristics of agency, denoting a human COLLECTIVE which also has a locational sense.

Table 4.15. The pairwise associations of verb-chain general morphological features on the one hand and syntactic arguments and their semantic and structural subtypes on the other hand which have been considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{2|1} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) and presented so that always $U_{2|1} \geq U_{1|2}$; $F_1 > F_2 \sim U_{2|1} > U_{1|2}$; $F_1 \subset F_2 \sim F_1$ is a logical subset of F_2 ; $F_1 \equiv F_2 \sim F_1$ is logically equivalent to F_2 so that $F_1 \subset F_2$ and $F_1 \supset F_2$; $F_1 \neq F_2 \sim F_1$ is logically multiply disjoint within a set of related features $\cup(F_1, \dots, F_n)$ so that $F_1 \not\subset \cup(F_2, \dots, F_n)$ and $(F_2, \dots, F_n) \not\subset F_1$; associations covered more generally by some other(s) or otherwise considered less informative in (parentheses); ‘?’ indicates (minor) inconsistency in underlying analysis scheme. Furthermore, correspondences arising among the various overlapping analysis schemes concerning verb-chains (SX_VCH.SEM_XXX, SX_XAUX, Z_ANL_XXX) or nonoccurrences of particular syntactic arguments (SX_XXX.SEM_NIL) mirroring a correspondence with a positive occurrence of the same argument (SX_XXX) are mostly ignored as noninformative.

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
SX_NAUX \equiv ?Z_ANL_NEG	0.863	0.856	314	310	298
Z_ANL_ACT \supset ?SX_AGE	0.69	0.622	2306	2537	2302
SX_AGE \neq ?Z_ANL_PASS	0.5	0.347	2537	457	11
Z_ANL_SING \subset ?SX_AGE	0.493	0.411	1962	2537	1961
(Z_ANL_SING \neq SX_AGE.SEM_NIL)	0.493	0.411	1962	867	1
(SX_AGE.SEM_INDIVIDUAL \subset Z_ANL_ACT)	0.469	0.461	2251	2306	2090
(Z_ANL_AFF \neq SX_NAUX)	0.421	0.233	2573	314	9
(SX_AGE.SEM_INDIVIDUAL \neq Z_ANL_PASS)	0.394	0.243	2251	457	6
Z_ANL_NEG $>$ SX_VCH.SEM_NONNECESSITY	0.354	0.068	310	36	33
(Z_ANL_SGPL3 $>$ SX_AGE)	0.322	0.266	1519	2537	1518
Z_ANL_THIRD \subset SX_AGE	0.322	0.266	1519	2537	1518
Z_ANL_NEG $>$ SX_VCH.SEM_SINENECESSITY	0.298	0.083	310	57	46
Z_ANL_COVERT $>$ SX_AGE.SEM_INDIVIDUAL	0.294	0.288	1218	2251	1214
Z_ANL_SING $>$ SX_AGE.SEM_INDIVIDUAL	0.293	0.275	1962	2251	1772
Z_ANL_OVERT $>$ SX_AGE.SEM_GROUP	0.287	0.115	1314	256	256
Z_ANL_NEG $>$ SX_VCH.SEM_IMPOSSIBILITY	0.262	0.098	310	83	60
Z_ANL_OVERT $>$ SX_AGE	0.262	0.223	1314	2537	1313
Z_ANL_SG3 \subset SX_AGE	0.246	0.212	1257	2537	1256
Z_ANL_AFF $>$ SX_VCH.SEM_NONNECESSITY	0.229	0.024	2573	36	1
Z_ANL_COVERT $>$ SX_AGE	0.227	0.197	1218	2537	1214
Z_ANL_PASS $>$ SX_LOC.SEM_GROUP	0.22	0.047	457	56	44
Z_ANL_OVERT $>$ SX_PAT.DIRECT_QUOTE	0.215	0.049	1314	120	119
SX_AGE.SEM_INDIVIDUAL \supset Z_ANL_SGPL12	0.201	0.174	2251	829	824

Finishing this scrutiny of pairwise associations with those concerning extra-linguistic features as at least one of the features, we can see in Table 4.16 that the strongest associations concern logically complementary or equivalent relationships such as between the two sub-corpora as well as some of the repetition-related features, or the restriction of one feature-category to only one of the two subcorpora, namely, quotations which occur only in the newspaper text. The only strong association which is not entirely logically predetermined is that between DIRECT QUOTES as PATIENT and the positioning of attributive structures with any of the THINK lexemes after such quotes; though this is predominantly the case the figures indicate that there is also a very small minority ($n=4$, 3.3%) of attributive structures which instead precede the quotation. With respect to other feature combinations involving quotations, we may see that THIRD PERSON SINGULAR in the verb-chain as well as OVERT manifestation of

the subject/AGENT are moderately associated with the attributive structures, which are characteristics one could expect in “reported speech”, while the letters-to-the-editor subsection (EXTRA_DE_hs95_MP) in the newspaper subcorpus does not contain such quoted passages at all.

The other relationships here indicate individual authors who have contributed only or predominantly to just one of the various subsections in either of the two subcorpora, for example, the one journalist at Helsingin Sanomat who has exclusively written articles published in the cultural section (EXTRA_DE_hs_95_KU) during the two-month period of the newspaper sampled into the research corpus. An interesting detail is that one of the contributors to the Internet newsgroup discussion (#966), who has firstly written only to the politics-related forum (EXTRA_DE_sfnet_politiikka), also accounts for a substantial portion of the occurrences of the studied THINK lexemes in the SECOND PERSON PLURAL or as part of verb chains with that particular feature. Finally, regarding potential repetition within individual texts, *ajatella* has the highest, though moderate, association of being followed by another THINK lexeme (though it may be any one of the four selected ones), which can be considered a natural consequence of this particular lexeme being by far the most frequent of the lot.

Table 4.16. The pairwise associations of extra-linguistic features both mutually and with other feature categories, which have been considered on the basis of the Uncertainty Coefficient to have at least a moderate relationship ($U_{2|1} > 0.2$), calculated asymmetrically as the association of the first mentioned feature F_1 with the second one F_2 ($U_{2|1}$) and vice versa ($U_{1|2}$) and presented so that always $U_{2|1} \geq U_{1|2}$; $F_1 > F_2 \sim U_{2|1} > U_{1|2}$; $F_1 \subset F_2 \sim F_1$ is a logical subset of F_2 ; $F_1 \equiv F_2 \sim F_1$ is logically equivalent to F_2 so that $F_1 \subset F_2$ and $F_1 \supset F_2$; $F_1 || F_2 \sim$ logical complementarity throughout the entire data so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$ and $\forall x(x \in F_1 \vee x \in F_2)$; $F_1 | F_2 \sim F_1$ is logically pairwise disjoint with F_2 so that $F_1 \not\subset F_2$ and $F_2 \not\subset F_1$; $F_1 \neq F_2 \sim F_1$ is logically multiply disjoint within a set of related features $\cup(F_1, \dots, F_n)$ so that $F_1 \not\subset \cup(F_2, \dots, F_n)$ and $(F_2, \dots, F_n) \not\subset F_1$; associations covered more generally by some other(s) or otherwise considered less informative in (parentheses); ‘?’ indicates (minor) inconsistency in underlying analysis scheme.

Feature pair	$U_{2 1}$	$U_{1 2}$	n_1	n_2	n_{common}
Z EXTRA_SRC sfnet Z EXTRA_SRC hs95	1	1	1654	1750	0
(Z PREV FIRST≡Z PREV NONE)	1	1	2641	2641	2641
SX PAT.DIRECT_QUOTE⊃Z POST_QUOTE	0.965	0.941	120	116	116
(Z EXTRA_SRC hs95⊃Z NON_QUOTE)	0.566	0.545	1750	1312	1312
(Z EXTRA_SRC sfnet≠Z NON_QUOTE)	0.566	0.545	1654	1312	0
(Z PREV NONE Z PREV REPEAT)	0.544	0.348	2641	364	0
(Z PREV FIRST Z PREV REPEAT)	0.544	0.348	2641	364	0
Z EXTRA_DE hs95_KU>Z EXTRA_AU hs95_kivi...	0.477	0.091	224	27	27
(Z SG3>Z POST_QUOTE)	0.41	0.144	509	116	112
Z EXTRA_DE politiikka>Z EXTRA_AU sfnet 966	0.365	0.083	626	77	77
Z EXTRA_DE politiikka>Z EXTRA_AU sfnet 948	0.299	0.032	626	30	30
Z EXTRA_DE ihmissuhteet>Z EXTRA_AU sfnet 92	0.258	0.046	1028	79	79
Z EXTRA_DE ihmissuhteet>Z EXTRA_AU sfnet 345	0.253	0.043	1028	73	73
Z PREV_REPEAT>Z PREV_ajatella	0.245	0.226	364	322	187
(Z EXTRA_AU sfnet 966>Z_PL2)	0.243	0.175	77	51	25
Z EXTRA_AU sfnet 966>Z_ANL_PL2	0.228	0.197	77	64	29
Z EXTRA_DE hs95_MP>Z EXTRA_AU hs95_pääte...	0.227	0.169	105	72	35
Z_NON_QUOTE>Z_EXTRA_DE_hs95_MP	0.219	0.045	1312	105	105
(Z EXTRA_SRC hs95⊃Z_QUOTE)	0.215	0.096	1750	318	318
(Z EXTRA_SRC sfnet≠Z_QUOTE)	0.215	0.096	1654	318	0
Z_ANL_SG3>Z_POST_QUOTE	0.213	0.048	1257	116	114
Z_ANL_OVERT>Z_POST_QUOTE	0.212	0.047	1314	116	115
Z EXTRA_DE ihmissuhteet>Z EXTRA_AU sfnet 855	0.208	0.016	1028	28	28
Z EXTRA_DE ihmissuhteet>Z EXTRA_AU sfnet 331	0.201	0.043	1028	99	92

4.2.2 Pairwise associations of grouped features

The grouped-feature analyses to be presented below have been motivated by the preceding both univariate and bivariate results. The first scrutiny presented in Table 4.17 concerns the relationships between the various major morphological categories of NON-FINITE verb forms, that is, INFINITIVES and PARTICIPLES, and the nominal cases. The association values measured both ways are relatively high, and the overall relationship is statistically significant; cell-wise, the preferences appear to confirm as well as supplement earlier results. So, the FIRST INFINITIVE remains neutral with respect to all cases, whereas the SECOND INFINITIVE has a positive preference for the INESSIVE in addition to the INSTRUCTIVE case noted earlier, while it disprefers all the other cases. Furthermore, the THIRD INFINITIVE exhibits an overall preference for only

the ILLATIVE case, associated with its obligatory government as an alternative infinitive form in verb-chains. Finally, the FOURTH INFINITIVE shows a preference for the NOMINATIVE, GENITIVE, and PARTITIVE cases, which is consistent with the alternative and presently ever more common interpretation of this form as rather a deverbal noun with *-minen* indicating ‘do → act of doing’.

With respect to the two participles, both exhibit a preference for the NOMINATIVE and PARTITIVE and a dispreference for the ILLATIVE and INSTRUCTIVE cases, the latter two which would appear to be of more use in conjunction with the infinitives in general. Moreover, while the PRESENT PARTICIPLE also prefers the GENITIVE case, the PAST PARTICIPLE instead prefers the TRANSITIVE and disprefers the INESSIVE case, with the GENITIVE remaining neutral this time. Among these, the NOMINATIVE can partially be traced to the use of both participles in compound tenses, and the TRANSITIVE with the ACCIDENTAL construction, while I would hypothesize on the basis of my native speaker intuition that the GENITIVE as well as the PARTITIVE cases might to a certain extent be linked to CLAUSE-EQUIVALENT forms based on the two participles, for example, *toivoisin asiaa*_{PATIENT+NOTION} *harkittavan*_{PARTICIPLE1+PASSIVE+GENITIVE} ‘I would hope the matter to be considered’ or *asiaa*_{PATIENT+NOTION} *harkittuani*_{PARTICIPLE2+PASSIVE+PARTICIPLE+FIRST} ‘having considered the matter’.

Table 4.17. Pairwise comparison of the subtypes of PARTICIPLES and INFINITIVES with morphological cases among the studied THINK lexemes; $P(df=36)=3.34e^{-272}$; $V_{\text{Cramér's}}=0.610$; $U_{\text{CASE|INFINITIVE/PARTICIPLE}}=0.538$; $U_{\text{INFINITIVE/PARTICIPLE|CASE}}=0.635$.
THINK.INFINITIVE PARTICIPLE vs CASE\$residual.pearson.std.sig

Features	Z NOM	Z GEN	Z PTV	Z TRA	Z INE	Z ILL	Z INS
Z INF1	0	0	0	0	0	0	0
Z INF2	-	-	-	-	+	-	+
Z INF3	-	-	-	-	-	+	-
Z INF4	+	+	+	-	-	-	-
Z PCP1	+	+	+	0	0	-	-
Z PCP2	+	0	+	+	-	-	-

Switching to the FINITE forms, I found that it would be interesting to study the relationships between 1) polarity and moods, 2) moods and person/number, and 3) person with number separated from each other. All of these relationships turn out to be statistically significant, evident in Tables 4.18-4.20. In particular, we can in Table 4.18 firstly note that AFFIRMATIVE polarity has a preference for the IMPERATIVE and a dispreference for the INDICATIVE mood, whereas the tables are turned in the case of NEGATIVE polarity, which has a positive preference for the INDICATIVE and a dispreference for the IMPERATIVE mood. In somewhat of a contrast, the CONDITIONAL mood is neutral with respect to both types of polarity. These results can be interpreted to indicate that NEGATION is relevant in recounting states-of-affairs, in this case specifically concerning non-existence, prototypical to the INDICATIVE mood, whereas commands and requests communicated with the IMPERATIVE are mostly positive exhortations instead of prohibitions.

Moving on to the relationships between person-number and mood presented in Tables 4.19, we can firstly see that the INDICATIVE has a preference for the FIRST PERSON SINGULAR, the THIRD PERSON PLURAL and the PASSIVE voice, whereas it exhibits a dispreference for both SECOND PERSON SINGULAR and PLURAL, with THIRD PERSON SINGULAR and FIRST PERSON PLURAL remaining neutral. In contrast, the IMPERATIVE is

almost a mirror image of the INDICATIVE mood, with a preference for both SECOND PERSON SINGULAR and SECOND PERSON PLURAL, with a dispreference for all the other person-number features except FIRST PERSON PLURAL which is neutral. For its part, the CONDITIONAL mood has a preference for only the THIRD PERSON SINGULAR, dispreferring FIRST and SECOND PERSON SINGULAR and THIRD PERSON PLURAL and the PASSIVE voice, with both FIRST and SECOND PERSON PLURAL staying neutral this time round. It is my judgement that the characteristic usage of the IMPERATIVE mood in issuing commands and requests to addressees in the immediate context is reflected in these results, as is the assignment of possible and tentative states-of-affairs associated with the CONDITIONAL mood to far-away contexts, detached from the “here and now”.

Table 4.18. Pairwise associations of the two POLARITY and three most common mood features among the verb-chains with the studied THINK lexemes; $P(df=2)=9.42e^{-05}$; $V_{\text{Cramér's}}=0.0812$; $U_{\text{MOOD|POLARITY}}=0.00887$; $U_{\text{POLARITY|MOOD}}=0.0138$.
THINK.POLARITY vs MOOD\$residual.pearson.std.sig

Features	Z ANL IND	Z ANL KOND	Z ANL IMP
Z ANL AFF	-	0	+
Z ANL NEG	+	0	-

Table 4.19. Pairwise associations of the person-number features as well as PASSIVE voice and the three most common moods in the verb-chains with the studied THINK lexemes; $P(df=12)=3.338e^{-258}$; $V_{\text{Cramér's}}=0.477$; $U_{\text{MOOD|PERSON+NUMBER}}=0.259$; $U_{\text{PERSON+NUMBER|MOOD}}=0.0890$.
THINK.PERSON NUMBER vs MOOD\$residual.pearson.std.sig

Features	Z ANL IND	Z ANL KOND	Z ANL IMP
Z ANL SG1	+	-	-
Z ANL SG2	-	-	+
Z ANL SG3	0	+	-
Z ANL PL1	0	0	0
Z ANL PL2	-	0	+
Z ANL PL3	+	-	-
Z ANL PASS	+	-	-

As a brief excursion to the separation of person and number as distinct atomic features instead of the person-number bundles they are typically analyzed as, we can next look at the mutual relationships between these two characteristics, presented in Table 4.20. It appears that the preferences are in this particular configuration restricted to the FIRST person, which has a positive preference for the SINGULAR and a dispreference for the PLURAL number. In contrast, the other two persons, both SECOND and THIRD, turn out to be neutral with respect to number. This result is most probably reflected in the association levels calculated in both directions, which are quite low as $U_{\text{PERSON|NUMBER}}=0.003$ and $U_{\text{NUMBER|PERSON}}=0.006$.

Table 4.20. Pairwise associations of the generalized PERSON and NUMBER features among the verb-chains with the studied THINK lexemes; $P(df=2)=0.00289$; $V_{\text{Cramér's}}=0.0706$; $U_{\text{NUMBER|PERSON}}=0.00584$; $U_{\text{PERSON|NUMBER}}=0.00295$.
THINK.PERSON vs NUMBER\$residual.pearson.std.sig

Features	Z ANL SING	Z ANL PLUR
Z ANL FIRST	+	-
Z ANL SECOND	0	0
Z ANL THIRD	0	0

Finally, the associations of the several subtypes of syntactic PATIENT and MANNER arguments noted above motivates a grouped-feature scrutiny of their subtypes in their entirety, presented in Table 4.21, using the most general classification scheme in the case of MANNER arguments. As can clearly be seen, substantial preferences or dispreferences are quite sparse, which is probably reflected in the relatively weak association levels, with PATIENT arguments bettering MANNER arguments as $U_{MANNER|PATIENT}=0.137$ but only $U_{PATIENT|MANNER}=0.0505$. Among the (positive) preferences we can find only ACTIVITIES and the closely related expressions of COMMUNICATION as PATIENT and the POSITIVE subtype of MANNER arguments, in addition to the nonoccurrence of both arguments having preferences with many of the other subtypes. The number of dispreferences is higher, with abstract NOTIONS as PATIENT evading GENERIC as well as AGREEMENT subtypes of MANNER, and ACTIVITIES as PATIENT shunning the same plus the lumped leftover category (OTHER1). Furthermore, among the structural subtypes of PATIENTS, INDIRECT QUESTIONS would appear to evade the GENERIC, FRAME, NEGATIVE, and AGREEMENT subtypes of MANNER, DIRECT QUOTES the GENERIC and POSITIVE subtypes of MANNER, and *että* ‘that’ clauses the GENERIC, FRAME, POSITIVE as well as AGREEMENT subtypes of MANNER. In contrast to these associations, human GROUPS, STATES, ATTRIBUTES, TIME, and INFINITIVES as subtypes of PATIENT arguments as well as the JOINT subtype of MANNER arguments exhibit no preferences or dispreferences at all with respect to the other argument type considered here.

Table 4.21. Pairwise comparison of the subtypes of PATIENT and MANNER arguments among the studied THINK lexemes; $P(df=105)=1.13e^{-101}$; $V_{Cramér's}=0.188$; $U_{MANNER|PATIENT}=0.137$; $U_{PATIENT|MANNER}=0.0505$.

THINK.PATIENT vs MANNER

Patient/ Manner	GEN...	FRA...	POS...	NEG...	AGR...	JOINT	OTHER1	NIL
INDIVIDUAL	0	0	0	0	0	0	0	+
GROUP	0	0	0	0	0	0	0	0
NOTION	-	0	0	0	-	0	0	0
STATE	0	0	0	0	0	0	0	0
ATTRIBUTE	0	0	0	0	0	0	0	0
TIME	0	0	0	0	0	0	0	0
ACTIVITY	-	0	+	0	-	0	-	+
EVENT	0	0	0	0	0	0	0	+
COMM...	0	0	+	0	0	0	0	0
INFINITIVE	0	0	0	0	0	0	0	0
PARTICIPLE	0	0	0	0	0	0	0	+
INDIRECT_Q...	-	-	0	-	-	0	0	+
DIRECT_Q...	-	0	-	0	0	0	0	+
<i>että</i> ‘that’ clause	-	-	-	0	-	0	0	+
OTHER	0	0	0	0	0	0	0	0
NIL	+	+	+	+	+	0	+	-

Here, we may recall that another argument combination concerning syntactic AGENTS and PATIENTS was already presented as an example demonstrating the paired grouped-feature analysis in Section 3.3.3. As that scrutiny did not include the miscellaneous rarer categories or the nonoccurrence of either argument type, we may verify those earlier results against the more comprehensive grouped-feature analysis of the subtypes of AGENTS and PATIENTS presented in Table 4.22 below. The number of

changes is quite small, probably due to the fact that the more frequent subtypes included in the earlier analysis already covered relatively comprehensively the occurrences of the studied THINK lexemes in the data. Here, a STATE or an EVENT as a PATIENT in conjunction with an INDIVIDUAL as an AGENT have turned into a dispreference from neutrality, while INDIRECT QUESTIONS and INFINITIVES as PATIENTS have become a preference in association with an INDIVIDUAL as an AGENT. Otherwise, this grouped-feature analysis links ATTRIBUTES as PATIENT with the lumped leftover category of AGENTS (OTHER) and the corresponding lump category of PATIENTS with the nonoccurrence of an AGENT. Furthermore, INFINITIVES, INDIRECT QUESTIONS, and DIRECT QUOTES prefer the occurrence of any type of AGENT, while NOTIONS and STATES prefer its nonoccurrence. In contrast, not having any type of PATIENT would be more non-preferable for a GROUP as an AGENT. Overall, knowing the subtype of PATIENT would explain clearly more of which is the subtype of AGENT than the other way around, as $U_{AGENT|PATIENT}=0.0720$ and $U_{PATIENT|AGENT}=0.0282$.

Table 4.22. Pairwise comparison of the subtypes of AGENT and PATIENT arguments among the studied THINK lexemes; $P(df=45)=6.08e^{-52}$; $V_{\text{Cramér's}}=0.198$; $U_{PATIENT|AGENT}=0.0282$; $U_{AGENT|PATIENT}=0.0720$.

THINK.AGENT vs PATIENT

Patient/Agent	INDIVIDUAL	GROUP	OTHER	NIL
INDIVIDUAL	0	0	0	0
GROUP	0	0	0	0
NOTION	-	+	0	+
STATE	-	0	0	+
ATTRIBUTE	0	0	+	0
TIME	0	0	0	0
ACTIVITY	-	+	0	+
EVENT	-	0	0	+
COMMUNICATION	0	0	0	0
INFINITIVE	+	0	0	-
PARTICIPLE	0	0	0	0
INDIRECT QUESTION	+	0	0	-
DIRECT QUOTE	+	-	0	-
<i>että</i> 'that' clause	+	-	0	-
OTHER	0	0	0	+
NIL	0	-	+	0

5 Multivariate analyses

5.1 Selection of variables

As was noted in Section 3.4.2, in order to ease comparative work now as well as later in the multivariate analyses I will fit and test several polytomous logistic regression models with varying degrees of complexity with the same data as is used with the final full model. These simpler models will include those containing 1) only node-specific morphological features, 2) verb-chain general morphological features as well as those node-specific features which are not subsumed by the verb-chain general ones, 3) syntactic argument types, *without* their semantic and structural subclassifications, 4) verb-chain general morphological features and non-subsumed node-specific morphological features together with syntactic argument types, the latter again *without* their subtypes, and 5) the aforementioned features and the most common semantic and structural classifications of AGENTS and PATIENTS, with the less frequent subtypes collapsed together whenever possible.

All of these less complex models will easily conform to the prescribed maximum of approximately 40 or so explanatory variables, determined by the overall number of instances of the least frequent outcome class, that is, in this study the lexeme *harkita*. For these models, the main remaining task in variable selection is to identify pairwise excessively strongly associated features and omit one of the two for each such pairing. For inclusion in such considerations, as was noted in Section 4.2.1 I have set the critical threshold value at $U_{21} > 0.5$ in either direction, though feature pairs with lower association values but nevertheless of overall general interest will also be shown and scrutinized. Despite this, in many individual cases I will just have to make a choice between two alternatives, which in principle are equally good. In such circumstances, I hope to be able to take into consideration and balance the overall makeup of the variable set. Furthermore, I will attempt to prefer more distinctive, “surprising” features over more prototypical, default cases.

For the node-specific morphological features, on the basis of the bivariate comparisons of which the relevant selection is presented in Table 5.1 below, this results in choosing the SECOND INFINITIVE over the INSTRUCTIVE case and the THIRD INFINITIVE over the ILLATIVE case, the latter which is also linked to the rejection of SINGULAR number (applicable for nominal-like NON-FINITE forms) in comparison to the THIRD INFINITIVE. On the node-specific level, I will retain some binary logically disjoint features which are not fully complementary over all occurrences of verbs, that is, both ACTIVE and PASSIVE voice as well as PAST and PRESENT tense. In contrast, the full complementarity between FINITE and NON-FINITE forms will obligatorily require the omission of at least one in this pair; however, as NON-FINITE forms are a general superset including all the INFINITIVES and PARTICIPLES, and the FINITE feature covers all verb forms with person and number, among others, I decided to rather omit both as too general and lacking added informative value from the linguistic perspective.

Table 5.1. The selection of morphological features on the basis of pairwise comparisons and other more general considerations; features excluded on the basis of the immediate comparison ~~struck through~~; features excluded on the basis of more general considerations ~~double struck through~~.

Feature pair	$U_{2 1}$	$U_{1 2}$
Z_NFIN Z_FIN	1	1
Z_INF2> Z_INS	0.866	0.75
Z_INF3> Z_ILL	0.748	0.676
(Z_SG>Z_INF3)	0.504	0.297

For the verb-chain general features, the selections of which are motivated in Table 5.2, the most unproblematic interpretation of the pairwise comparisons is to choose the combined SECOND person, more distinctive in my opinion, over the IMPERATIVE mood. As was noted earlier, all three moods exhibit a substantial level of mutual association, that is, multicollinearity, with $U_{Z_ANL_KOND|Z_ANL_IND}=0.378$ and $U_{Z_ANL_IMP|Z_ANL_IND}=0.309$, but these values are too low by themselves to warrant the exclusion of either of the two remaining ones, namely, the INDICATIVE and the CONDITIONAL. One could consider leaving out the most common of the moods, INDICATIVE, due to its relatively strong disjoint association with CLAUSE-EQUIVALENT forms in general, but as these features concern two distinct verb construction types as well as disparate levels of analysis and they are not fully complementary, even if counting in the two less frequent moods, I decided in the end to retain both features. The omission of ACTIVE voice is motivated on the basis that it is a logical superset of all FINITE forms with person and number, in which case the latter are in my view again linguistically more informative together with its counterpart, PASSIVE voice.

Likewise, such greater distinctiveness also accounts for my selection of NEGATION over AFFIRMATIVE polarity, as well as the COVERT expression of subjects/AGENTS over their OVERT manifestation, as the former feature is associated with more personal as well as situationally variable choice on the behalf of the writer/speaker. The same aspect also applies in choosing PLURAL over SINGULAR number (N.B. here originating from the person-number features in FINITE verb forms in a verb-chain or the corresponding possessive suffixes in conjunction with NON-FINITE CLAUSE-EQUIVALENT forms but *not* the nominal inflection of other NON-FINITE forms). This is further motivated by the multiple overlappings of the FIRST vs. SECOND vs. THIRD person, SINGULAR vs. PLURAL number, and PASSIVE voice in which at least one feature, in this case SINGULAR number, is always fully redundant and deducible from the values of the rest. Moreover, I had already opted above to include the three person features in viewing them as more informative. When combined with the node-specific morphological features, the aforementioned selections, or rather the rejections, on the verb-chain general level entail that the corresponding subsumed node-specific features are naturally also omitted, thus concerning IMPERATIVE mood, ACTIVE voice as well as all the specific person-number features (i.e., FIRST PERSON SINGULAR, SECOND PERSON SINGULAR, and so forth), in addition to the node-specific rejections already covered above.

Table 5.2. The selection of verb-chain general morphological features on the basis of pairwise comparisons and other more general considerations; features excluded on the basis of the immediate comparison ~~struck through~~; features excluded on the basis of other non-immediate comparisons or more general considerations ~~double struck through~~.

Feature pair	$U_{2 1}$	$U_{1 2}$
Z ANL SECOND> Z ANL IMP	0.585	0.342
Z ANL IND≠Z PHR CLAUSE	0.516	0.362
Z ANL SING<?Z ANL ACT	0.509	0.47
Z ANL AFF Z ANL NEG	0.471	0.259
Z ANL ACT Z ANL PASS	0.445	0.279
Z ANL IND≠Z ANL KOND	0.378	0.174
Z ANL OVERT Z ANL COVERT	0.36	0.352
Z ANL SING≠Z ANL PASS	0.329	0.191
Z ANL IND≠ Z ANL IMP	0.309	0.092
Z ANL SING Z ANL PLUR	0.304	0.158

The syntactic arguments alone have only one essential association, concerning that of CO-ORDINATED VERBS with CO-ORDINATED CONJUNCTIONS (Table 5.3). In this particular case, the CO-ORDINATED VERBS contain more variation, as is exhibited in their semantic subtypes, and are consequently selected, even more so as the associated CO-ORDINATED CONJUNCTIONS may sometimes be omitted. When considering the combination of node-specific and verb-chain general morphological features with syntactic arguments and their subtypes, the pairwise associations in Table 5.4 yield more data on the basis of which to select or omit variables. The most straight-forward cases here are those concerning individual morphological features as exponents of more general semantic characterizations of the verb-chains they pertain to, in which the selection of ACCIDENTAL verb-chains over TRANSLATIVE case, as well as EXTERNAL cause over both the ILLATIVE case and the THIRD INFINITIVE should need no further motivation. The same applies also to the selection of NEGATION as a verb-chain general feature over the practically equivalent occurrence of an explicit negated auxiliary (SX_NAUX) in the verb-chain.

We may next note that the prior rejection of both ACTIVE voice and SINGULAR number is further motivated by their strong association with the syntactic AGENT in general, and its INDIVIDUAL subtype in particular. One might also contemplate omitting verb-chain general PASSIVE voice in conjunction with the syntactic arguments due to its strong association with the AGENT, as $U_{Z_ANL_PASS|SX_AGE}=0.5$, though the relationship is not fully complementary as 421 (12.4%) occurrences of the studied THINK lexemes have neither an explicit specific AGENT nor an implicitly expressed, unspecified human one indicated by the PASSIVE voice. Like the three moods discussed above, this is also an example of multicollinearity, the empirical fact that the more variables one uses the more they typically are interrelated in one way or another. Another case in the same vein which is exemplified in Table 5.4 is that modality of some sort, or alternatively, the absence of any subtype of modality, has a relatively strong relationship with the FIRST INFINITIVE, with $U_{Z_INFI|SX_VCH.SEM_NILMODALITY|Z_INFI}=0.394$ and $U_{SX_VCH.SEM_NILMODALITY|Z_INFI}=0.359$, and to a lesser extent with the more general NON-FINITE forms as a whole, with $U_{SX_VCH.SEM_NILMODALITY|Z_NFIN}=0.286$ and $U_{Z_NFIN|SX_VCH.SEM_NILMODALITY}=0.234$.

One can on the basis of this motivate the exclusion of structural (node-specific) morphological features, that is, the different types of INFINITIVES and PARTICIPLES as

well as nominal cases, altogether from the full model incorporating the overall semantic classifications of the verb chains as well as the verb-chain general morphological features, since the particular morphological forms of the various components constituting the entire verb-chain are largely determined by idiosyncratic though mostly regular grammatical rules of Finnish, and lack much semantic content on their own. For instance, an individual lexeme as an auxiliary verb determines which one of the two common alternative INFINITIVES, the FIRST or the THIRD, should be used. Nevertheless, while I will exclude the specific type of infinitive or participle from the full model, I will retain the general feature representing the usage of THINK lexemes as the node verb of a CLAUSE-EQUIVALENT construction (Z_PHR_CLAUSE). Among the infinitives this feature will in general apply for the SECOND (associated with the INSTRUCTIVE case as noted above) and FOURTH INFINITIVES, which both have a small likelihood of occurring as part of a FINITE verb-chain in comparison to the FIRST and THIRD INFINITIVES, manifest in the association values with $U_{Z_INF2|Z_PHR_CLAUSE}=0.392$ and $U_{Z_INF4|Z_PHR_CLAUSE}=0.306$.⁸⁵ However, considerably more linguistically informative in my view would be the semantic classes of these CLAUSE-EQUIVALENTS, even according to the conventional Finnish grammar as the participial, temporal, agent, and purposive constructions etc. (Karlsson 1983, 2008), had such classifications been undertaken for the data at hand and had there been more leeway for explanatory variables.

Table 5.3. The selection of syntactic argument types on the basis of pairwise comparisons; features excluded on the basis of the immediate comparison ~~struck through~~.

Feature pair	U _{2 1}	U _{1 2}
SX CV<?SX CC	0.837	0.761

Table 5.4. The selection of node-specific verb-chain general morphological features in combination with the syntactic arguments and their subtypes on the basis of pairwise comparisons and other more general considerations; features excluded on the basis of the immediate comparison ~~struck through~~; features excluded on the basis of other non-immediate comparisons or more general considerations ~~double struck through~~.

Feature pair	U _{2 1}	U _{1 2}
Z TRA>SX VCH.SEM ACCIDENTAL	0.807	0.684
SX NAUX=?Z ANL NEG	0.863	0.856
Z ANL ACT>?SX AGE	0.69	0.622
SX AGE=?Z ANL PASS	0.5	0.347
Z ANL SING<?SX AGE	0.493	0.411
(SX AGE.SEM INDIVIDUAL<Z ANL ACT)	0.469	0.461
Z ILL>SX VCH.SEM EXTERNAL	0.439	0.176
Z INF3>SX VCH.SEM EXTERNAL	0.421	0.153
SX VCH.SEM NILMODALITY>Z INF1	0.394	0.359
(SX AGE.SEM INDIVIDUAL≠Z ANL PASS)	0.394	0.243

With respect to the selection of extra-linguistic features presented in Table 5.5, among the two complementary sources I have decided to select the one representing Internet newsgroups discussion (Z_EXTRA_SRC_sfnet), since as a more informal register it

⁸⁵ The only non-clause-equivalent verb-chains that I can think of which would contain the SECOND INFINITIVE is the quite idiosyncratic construction *mikä on ollessa/ajatellessa* ‘what is there [then to be annoyed] in being/thinking’ studied by Kotilainen (2007b); a similar and equally rare example for the FOURTH INFINITIVE as a component of a verb-chain would be its use in the archaic neccessive construction *minun on ajattelemisen* ‘it is my part/task to think’ discussed in Appendix C.

can be neatly coupled with quotations within newspaper text (Z_QUOTE), which for the most part represent spoken language. As the attributive constructions using a THINK lexeme almost categorically follow the quotation that they refer to, the feature referring to the position of the attributions becomes redundant (Z_POST_QUOTE), as is also the case with “normal” unquoted text, pertaining categorically only to the newspaper subcorpus.

What comes to features concerning repetition, I have in Table 5.5 rejected not having a preceding THINK lexeme within a text (Z_PREV_NONE) as the full equivalent of being the first THINK lexeme within such a text (Z_PREV_FIRST). Nevertheless, I have overall come to the conclusion to exclude them all from these multivariate considerations, since the prior univariate analyses have demonstrated that they rank feature-wise as well as lexeme-wise among the second-lowest in terms of explanatory power in comparison to the other feature categories, with a mean $U_{L|F}=0.002$ and $U_{F|L}=0.018$, exceeding only author identities in this respect. This decision is further facilitated in that “firstness” and repetition of a particular previous THINK lexeme within the same text have a relatively strong mutual association, which is quite understandable as the number of THINK lexemes per any of the relatively short coherent texts in the research corpus is overall low, with only 549 (20.8%) of the texts containing *more* than one THINK lexeme (out of the altogether 2641 distinct texts in the research corpus with at least one THINK lexeme).

Table 5.5. The selection of (mainly) extra-linguistic features on the basis of pairwise comparisons and other more general considerations; features excluded on the basis of the immediate comparison ~~struck through~~; features excluded on the basis of other non-immediate comparisons or more general considerations ~~double struck through~~.

Feature pair	$U_{2 1}$	$U_{1 2}$
Z EXTRA_SRC sfnet Z EXTRA_SRC hs95	1	1
(Z PREV_FIRST=Z PREV_NONE)	1	1
SX PAT.DIRECT QUOTE →Z POST_QUOTE	0.965	0.941
(Z EXTRA_SRC hs95→Z NON_QUOTE)	0.566	0.545
(Z EXTRA_SRC sfnet ≠Z NON_QUOTE)	0.566	0.545
(Z PREV_NONE Z PREV_REPEAT)	0.544	0.348
(Z PREV_FIRST Z PREV_REPEAT)	0.544	0.348
(Z EXTRA_SRC hs95→Z QUOTE)	0.215	0.096
(Z EXTRA_SRC sfnet ≠Z QUOTE)	0.215	0.096

However, the overall number of contextual variables evident in the data and potentially incorporable in the full model even after the initial pruning and selection of the most general and higher-level features (i.e., syntactic arguments and in particular their semantic and structural subtypes) over more specific and purely morphosyntactic ones at the end of Section 4.1.1, as well as subsequent to the selections and rejections on the basis of the pairwise comparisons conducted immediately above, is still closer to one hundred, if one considers the finer-grained levels of analysis applied for some syntactic argument types. This clearly exceeds the limits recommended for logistic regression analysis with the available data set. Therefore, I will still have to undertake some further drastic reductions in comparison to the full range of intricacy applicable in the univariate analyses and apparent in the results presented therein. In this, knowledge of the subject matter is preferred (see, e.g., Harrell 2001: 66), but as this is an exploratory study with little prior research

existing which concerns specifically the studied THINK lexemes, I will have to resort to my general professional understanding of the linguistic analysis scheme I have applied. Moreover, the fact that some frequent feature does not have a significant distribution is not a very good motivation for its rejection, since that might distort and inflate the effects estimated for the remaining variables (Harrell 2001: 60).

Subsequently, in the proper full model I will include only verb-chain general features and their most general semantic classifications but no node-specific ones, and with respect to syntactic arguments I will incorporate their semantic and structural subtypes at the highest level of granularity *and* only in the case of the most frequent ones, i.e., PATIENT, AGENT, MANNER, and TIME-POSITION, while the rarer will be incorporated simply and uniquely as a syntactic argument type, i.e., META-comments (clause-adverbials), LOCATION, CO-ORDINATED VERBS, DURATION, FREQUENCY, QUANTITY, SOURCE, and GOAL, even if subtypes were available for them. In the case of some syntactic arguments, specifically SOURCE, QUANTITY, and CO-ORDINATED VERBS, and to a somewhat lesser extent also DURATION, their subtypes do in fact happen to have quite similar preference patterns, so the effect of this pruning is probably not that detrimental. However, for other arguments such as LOCATION and FREQUENCY with clearly distinctive subtypes this may lead to substantial loss of linguistically interesting information. The same concern applies also to some contrasted subtypes of verb-chain modality, but as not all of these are sufficiently frequent, for example, FUTILITY among the three subtypes of NECESSITY and STOP among TEMPORAL ones, I will stick to the most general classes in their case, too.

Furthermore, I will collapse a subset of the already quite numerous semantic subtypes of PATIENT into two more general ones, namely, HUMAN referents (including both INDIVIDUALS and GROUPS) and ABSTRACTIONS (including abstract NOTIONS, STATES, ATTRIBUTES, and TIME), as these appeared convergent in terms of their preference patterns in the univariate scrutiny and also form linguistically motivatable supersets. Likewise, I will also merge PURPOSE with REASON among the syntactic argument types without semantic subtypings. In the end, this leaves us still with altogether 46 explanatory variables in the final model proper. This variable set is somewhat smaller and different in feature type in comparison to the one used in the preliminary version of the results of this study presented in Arppe (2007), with altogether 59 explanatory variables. The main difference is in the selection of syntactic argument features alone instead of their sufficiently frequent semantic subtypes in the case of the rarer arguments. This choice should increase the overall proportion of the studied THINK lexemes covered by the features and the number of features associated with each lexeme, though it may in some cases lead to a loss of semantic precision, i.e., FREQUENCY as an argument type vs. its OFTEN and AGAIN subtypes. Furthermore, I have also opted for the more general features of modality, i.e., POSSIBILITY and NECESSITY, instead of their opposed subtypes, while now including the associated EXTERNAL cause. Moreover, I have excluded IMPERATIVE voice on the basis of further considerations of the pairwise associations.

However, since I am intrigued by what results might be produced with the entire variable set containing all the semantic and structural subtypes of the syntactic arguments satisfying the minimum frequency requirement and how they might compare with the prior univariate results, because the only real cost is computational, I will also try out such an extended model, even at the risk of not setting the best

example in the methodological sense. This extended model conforms more in its size and composition to the one used in Arppe (2007), though I have now also included some of the lumped subtypes for the less frequent syntactic arguments in order to increase overall the number of features associated with each lexeme in the data set. Nevertheless, variable clustering using statistical techniques such as Principal Components Analysis (PCA) or Hierarchical Cluster Analysis (HCA) would present the next step forward to further prune down the number of variables (see, e.g., Harrell 2001: 66-67),⁸⁶ but I have ruled out the use of such methods in this study as the number of various analysis stages is already extensive.

In general with respect to rarer subtypes presented in conjunction with the univariate analysis, these can sometimes be lumped together and coherently reinterpreted (by definition on the basis of the subtypes) in a linguistically meaningful way, but sometimes this cannot be achieved. In the former case, I will include such collapsed classes alongside the more frequent ones in the extended model, denoting the countable but NON-OFTEN subtype of FREQUENCY, and DURATION with a FIXED TEMPORAL REFERENT demarcating either or both ends for the time-period in question. In the latter case, even though such potentially unifying characterizations may be emergent they might not necessarily be uniformly applicable, in which case preference patterns may result simply as a product of chance. Thus, I will rather exclude such collapsed categories to be on the safe side, which concerns the rarer subtypes of AGENT, and PATIENT as well as LOCATION and MANNER. The variable sets for all the different models presented and discussed here above are presented in their entirety in Appendix R, whereas their general composition is summarized in Table 5.6 below. In addition, sets containing the extra-linguistic variables both alone and with the proper and extended full models have been included.

⁸⁶ This has been suggested to me by Dirk Speelman and Kris Heylen on separate occasions.

Table 5.6. Composition of the various features sets to be covered in the multivariate analyses as explanatory variables.

Model index	Feature set composition	Overall number of features
I	Only node-specific morphological features	26
II	Verb-chain general morphological features (10) as well as those node-specific features which are not subsumed by the verb-chain general ones (17)	27
III	Syntactic argument types, <i>without</i> their semantic and structural classifications	18
IV	Verb-chain general morphological features (10) and non-subsumed node-specific morphological features (17) together with syntactic argument types (17), the latter again <i>without</i> their subtypes	44
V	Verb-chain general features (10), the most common semantic classifications of AGENTS and PATIENTS with their less frequent subtypes collapsed together (12), and the other syntactic argument types alone <i>without</i> their subtypes (15)	37
VI	Proper full model with verb-chain general morphological features (10) and their semantic classifications (6) together with syntactic argument types alone (10) or their selected or collapsed subtypes (20)	46
VII	Proper full model with verb-chain general morphological features (10) and their semantic classifications (6) together with syntactic argument types alone (10) or their subtypes (20) as well as extra-linguistic features (2)	48
VIII	Extended full model with verb-chain general morphological features (10) and their semantic classifications (9) together with syntactic argument types (5) and all their subtypes exceeding the minimum frequency threshold (38)	62
IX	Extended full model with verb-chain general morphological features (10) and their semantic classifications (9) together with syntactic argument types (5) and all their subtypes exceeding the minimum frequency threshold (38) as well as extra-linguistic features (2)	64
X	Extralinguistic features alone (2)	2
XI	Syntactic argument types alone (10) or their selected or collapsed subtypes (20), together with semantic classifications of verb chains (6) but <i>without</i> any node-specific or verb-chain general morphological features	36

Having now fixed the variable sets we can at this stage evaluate to what extent the selected explanatory variables in particular are associated with each other. As we can see in Figure 5.1, the mean pairwise associations (calculated using the asymmetric Uncertainty Coefficient) among the variables selected in the proper full model are quite low, ranging $U_{2|1}=0.001-0.050$, as is the case also with the bulk of the maximal associations. Nevertheless, as can be summed on the basis of the preceding exposition some level of not insubstantial mutual interrelationship remains, with the maximum $U_{2|1}=0.516$ between verb-chain general INDICATIVE mood (Z_ANL_IND) and CLAUSE-EQUIVALENT usage the node-verb (Z_ANL_PHRASE), which are, however, in a disjoint relationship. However, the association levels drop then sharply, so that for 90 percent $U_{2|1}\leq 0.268$, that is, these are moderate relationships at best.

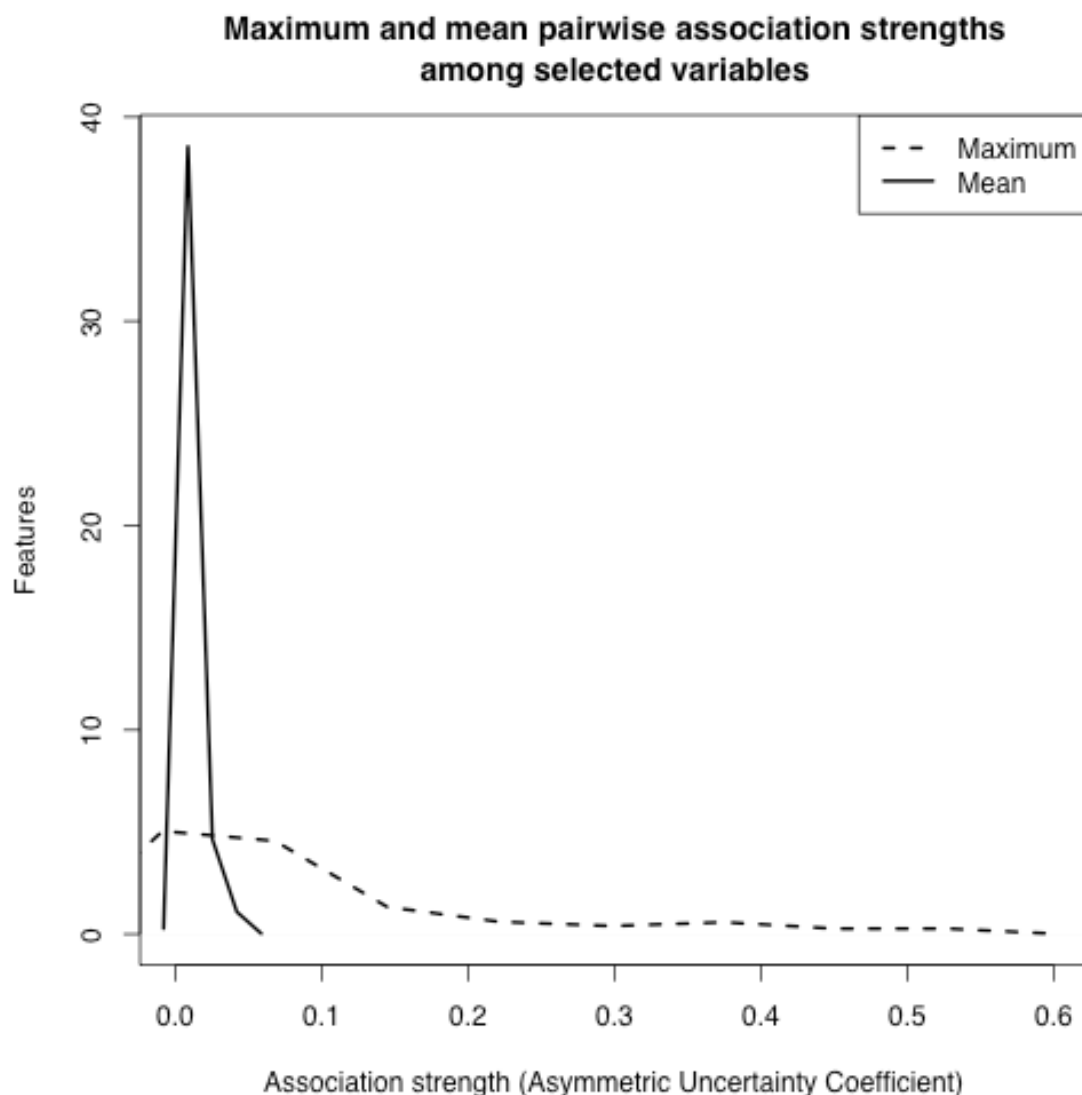


Figure 5.1. Maximum and mean pairwise association strengths, calculated using the asymmetric Uncertainty Coefficient ($U_{2|1}$), among the explanatory variables selected to the proper full model (VI).

We can also ascertain to what extent these selected variables are spread out in the data. Only six of the features do not have at least a single occurrence for all four of the selected THINK lexemes, but just one of these features does not have such occurrences for three (or more) of the lexemes, namely, GENERIC types of MANNER appear neither with *pohtia* nor with *harkita*. Furthermore, in the research corpus DIRECT QUOTES as PATIENT as well as the AGREEMENT subtype of MANNER and ACCIDENTAL verb-chains elude *harkita*, while INFINITIVES and PARTICIPLES as PATIENT shun *miettiinä*. The fact that *harkita* is somewhat prominent among these nonoccurrences can probably be attributed to its lowest frequency relative to the other THINK lexemes. From Figure 5.2 we can further deduce from the tips of the curves representing the relative proportions of the studied THINK lexemes with respect to the selected explanatory variables that there are only a few feature variables for which the most common lexeme gobbles up all occurrences with the feature in question. Rather, the mean proportion of the feature-wise most common lexemes is around one-half, and even the least frequent lexemes per each feature have as their mode proportion

approximately one-tenth. It is these proportions that a polytomous logistic regression model attempts to mimic, so once we have applied this particular statistical technique in the following Sections 5.2-5.3 we may evaluate in Section 5.5 how well the resultant models conform to this initial state-of-affairs represented in Figure 5.2. However, logistic regression modeling contrasts the occurrences of each feature with all the others present in the overall set of contexts, whereas the proportions in Figure 5.2 concern only singular features with no consideration for joint effects.

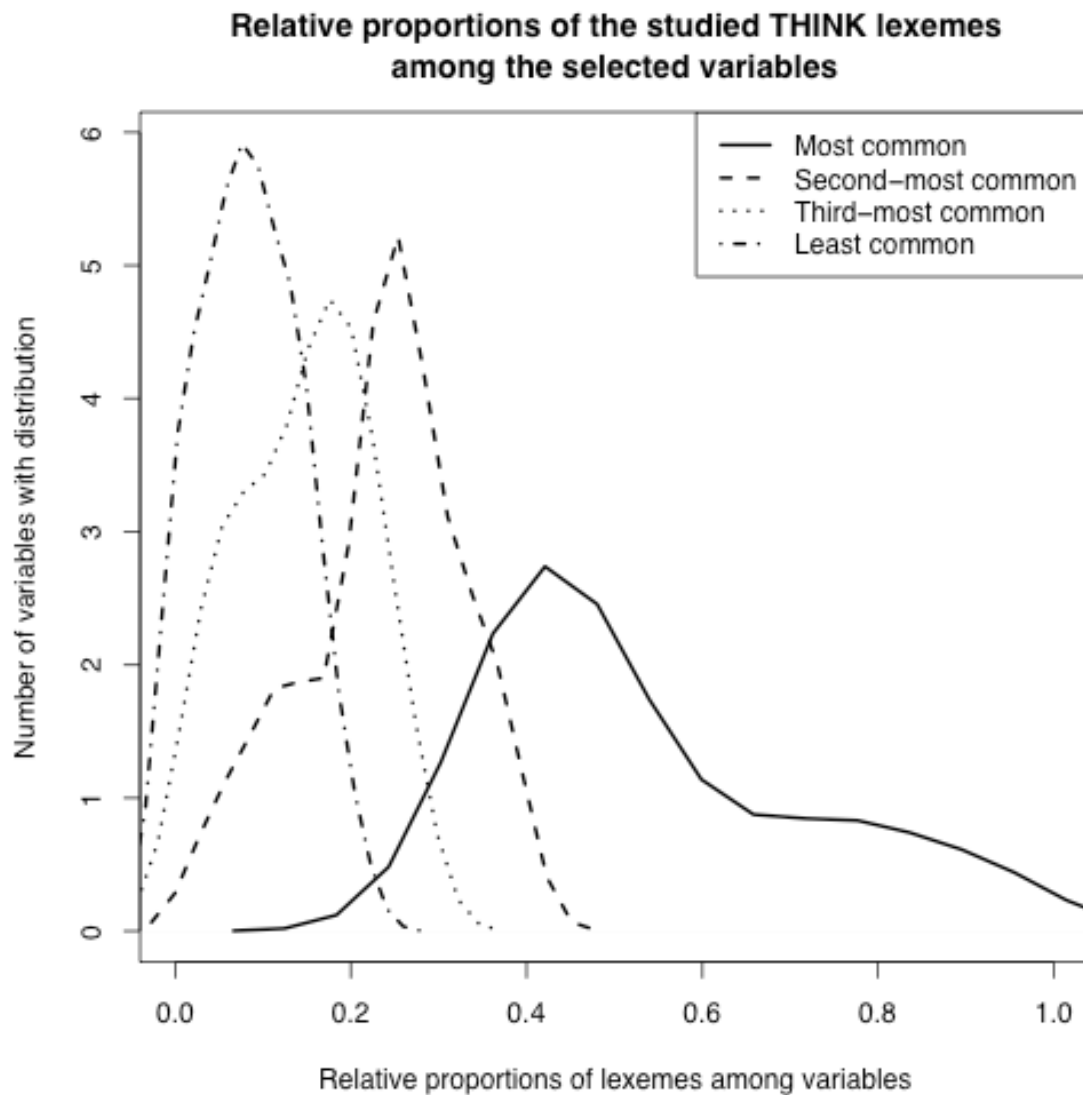


Figure 5.2. Relative proportions of the studied THINK lexemes among the selected explanatory variables in the proper full model (VI).

Finally, we can examine the number of occurrences of the feature variables selected in the proper full model per each of the studied THINK lexemes in the research corpus. As can be seen in Figure 5.3, fortunately all occurrences of the lexemes are associated with at least one of the selected explanatory variables. The maximum number of features occurring in a singular context is 11, but this applies in only three cases. The intermediate quartiles of feature occurrences are 4 (25%), 5 (50%), 6 (75%), which is also apparent in the Figure. What this also entails is that on the average only a

relatively small proportion of all the selected variables apply for any individual occurrence of the studied THINK lexemes.

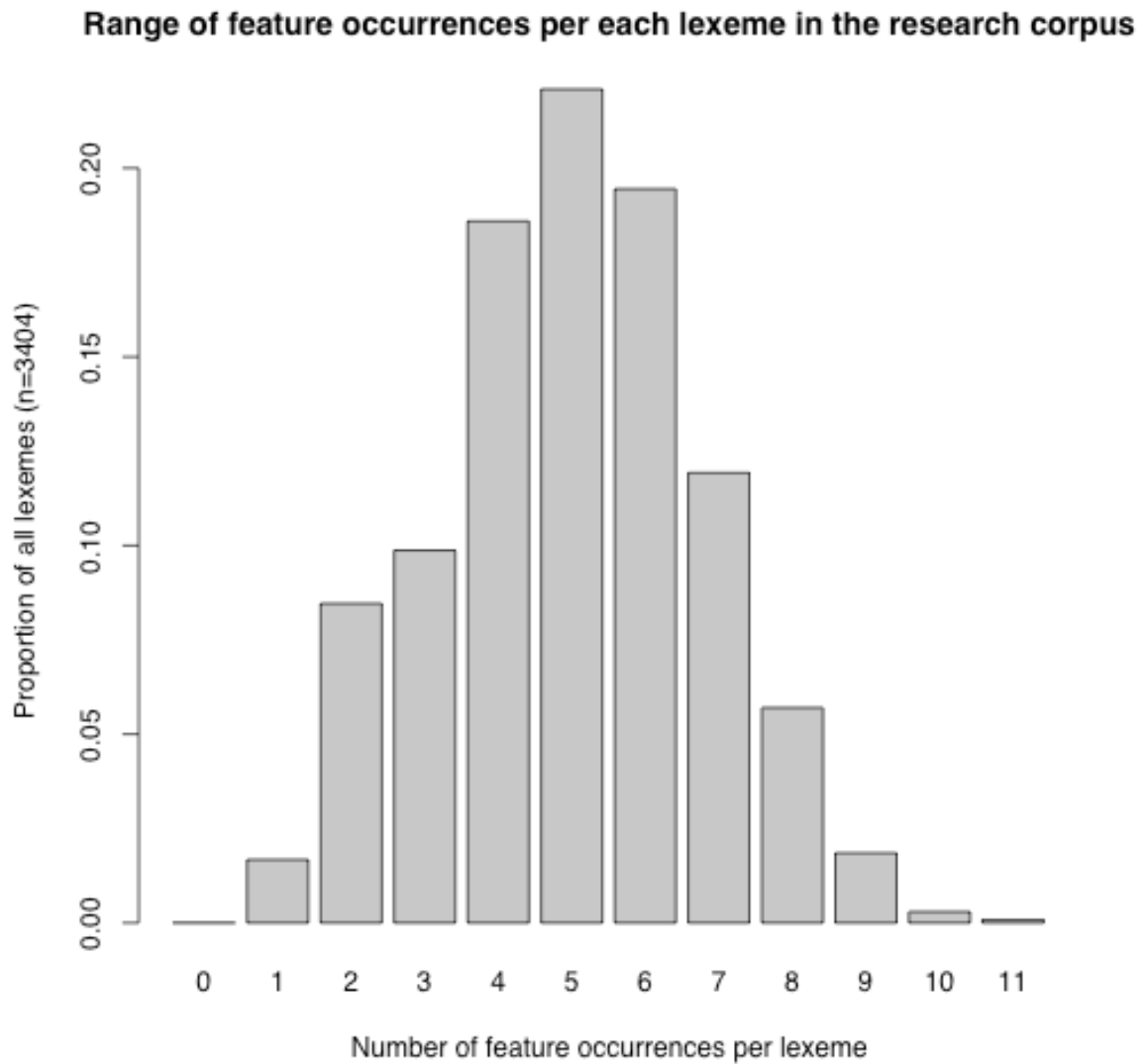


Figure 5.3. Range of the number of occurrences of the selected features in the proper full model (VI) in conjunction with individual contexts of the studied THINK lexemes in the research corpus, as a proportion of their overall frequency ($n=3404$).

5.2 Comparisons of the descriptive and predictive capabilities of the different heuristics and models

5.2.1 Comparing the various heuristics with respect to the full model

In comparing the performance of the various different heuristics presented in Sections 3.4.1 and 3.4.3 for implementing polytomous logistic regression, I will use as the reference model the proper full one (VI) described in the previous Section 5.1. The overall results of fitting this particular model with all these heuristics and testing their descriptive conformance as well as predictive capabilities with the same data set in its entirety are presented in Table 5.7 below. In order to assess the *process* of fitting the model using the various heuristics, I have in addition to the simple one-time fit and testing, using the entire data set in both circumstances, also validated the results with 1000-fold simple bootstraps, for which the results are shown in Table 5.8 further down. This choice differs from random grouped cross-validation used by Bresnan et al. (2007), which had a focus on the generalizability of a particular model on new, unseen data. Though bootstrapping can be more biased than cross-validation in favor of the model, the former has considerably less variance, that is, is more consistent, than cross-validation, when the entire validation process is repeated (Harrell 2001: 81, 90-96).

As can be seen in Table 5.7, all of the performance values for each of the considered heuristics do not differ to any substantial degree. Nevertheless, the simultaneously fitted proper multinomial appears best overall, followed in close procession first by ensembles of nested dichotomies (END), and then the pairwise and the one-vs-rest heuristics. Interestingly, among the specific individual nested dichotomies that the END heuristic builds upon, not even the two very best partitions $\{A, \{H, \{M, P\}\}\}$ or $\{P, \{A, \{M, H\}\}\}$ reach the same level as their aggregate, while the worst-rated partition $\{\{A, P\}, \{M, H\}\}$ is not substantially much lower.⁸⁷ This can probably be explained by the fact that as the END model is aggregated for each *instance* of the four outcomes, in this case the studied THINK lexemes, when even the best individual partition might considerably underperform systematically in some particular contexts, this can be offset by the better performance of some other partitions in such contexts, even though overall, and thus in a larger proportion of contexts, these other partitions perform worse. In other words, such smoothing which the END heuristic achieves is reflected in it performing overall better than any of its constituent partitions, and consequently these results support the advocacy of the END heuristic by Frank and Kramer (2004).

⁸⁷ On the basis of the existing lexicographical descriptions and my own intuition as a native speaker of Finnish, my best guess for the optimal partition would have been $\{A, \{\{M, P\}, H\}\}$, which the data through its analysis also raises to the top. The worst nested dichotomy is in my opinion also logical in that it first groups together two odd couples, namely $\{A, P\}$ and $\{M, H\}$. However, as the validated results in Table 5.8 indicate, the relative ranking of the various nesting partitions appears fluid, and the performance differences between the best and worst partitions remain on the average minimal.

Table 5.7. Performance of the various heuristics for polytomous logistic regression in both fitting and predicting with the proper full model (VI) using the original data in its entirety ($n=3404$).

Heuristic	R_L^2	Recall (%)	$\lambda_{prediction}$	$\tau_{classification}$
one-vs-rest	0.313	64.60	0.370	0.490
pairwise	NA	64.63	0.370	0.490
(simultaneous) multinomial	0.316	64.89	0.375	0.494
ensemble of nested dichotomies (END)	0.315	64.78	0.373	0.493
“best” nested dichotomies: {A, {H, {M, P}}} and {P, {A, {M, H}}}	NA	64.66	NA	NA
“worst” nested dichotomy: {{A, P}, {M, H}}	NA	63.66	NA	NA

The validation performance figures presented in Table 5.8 using simple bootstrap resamples fall slightly below those achieved by training with the entire data set, which can be attributed to each resample containing somewhat less of the overall variation apparent in the entire data, and consequently relatively fewer exemplars of the rarer, possibly exceptional or uncommon usages and contexts. The order of the heuristics in terms of their performance (*Recall*) is now slightly different, with the simultaneous multinomial falling to the lowest rank though practically similar to both the one-vs-rest and pairwise heuristics, while END is more distinctly separate from the rest at the apex. Nevertheless, the differences per each statistic and each technique are both minimal and consistent, thus suggesting that they do not essentially diverge as to their performance, at least with the particular linguistic phenomenon and selected variable set.

This conclusion is supported partially by the comparison of the absolute numbers of correctly predicted lexemes for the altogether 3404 outcomes, for which two-tailed *t*-tests between the performance of any pairing among the one-vs-rest, pairwise, and simultaneous multinomial heuristics indicate that their differences are not statistically significant (one-vs-rest vs. pairwise: $t=-0.763$, $df=1969.63$, $P=0.446$; pairwise vs. simultaneous multinomial: $t=-0.1895$, $df=1995.16$, $P=0.850$; one-vs-rest vs. simultaneous multinomial: $t=-0.980$, $df=1984.44$, $P=0.327$). However, the END heuristic appears to keep a distance which is statistically significant (END vs. one-vs-rest: $t=-5.288$, $df=1997.89$, $P=1.38e^{-07}$), albeit in absolute terms it remains still quite close to the other heuristics, too. Furthermore, a corresponding comparison of the R_L^2 figures, which represent the overall adherence of the estimated probabilities with the actual original outcomes, does show that the differences between all of these heuristics are in this respect nevertheless mutually significant (i.e., between one-vs-rest and simultaneous multinomial $t=-13.70$, $df=1816.46$, $P<2.2e^{-16}$, and between simultaneous multinomial and END $t=6.969$, $df=1928.09$, $P=4.355e^{-12}$).

Table 5.8. Validation of the performance of the various heuristics for polytomous logistic regression with respect to fitting and predicting with the proper full model (VI) using a 1000-fold simple bootstrap with the original data in its entirety ($n=3404$); Confidence Intervals (in parentheses) calculated using the percentile method.

Heuristic	R_L^2	Recall (%)	$\lambda_{prediction}$	$\tau_{classification}$
one-vs-rest	0.287 (0.264, 0.300)	63.80 (63.07, 64.51)	0.355 (0.343, 0.368)	0.479 (0.468, 0.489)
pairwise	NA	63.79 (62.87, 64.57)	0.355 (0.339, 0.369)	0.478 (0.465, 0.490)
(simultaneous) multinomial	0.292 (0.276, 0.302)	63.78 (62.96, 64.51)	0.355 (0.340, 0.368)	0.478 (0.466, 0.489)
ensemble of nested dichotomies (END)	0.294 (0.277, 0.305)	63.89 (63.10, 64.63)	0.357 (0.343, 0.370)	0.480 (0.468, 0.490)
“best” nested dichotomy: {A, {H, {M, P}}}	NA	63.65 (62.87, 64.37)	NA	NA
“worst” nested dichotomy: {A, {P, {M, H}}}	NA	63.01 (61.93, 63.84)	NA	NA

Another aspect in the performance of the different heuristics is to what extent they predict the same lexemes or not, shown in Table 5.9 below. As can be seen, the four considered heuristics are very convergent when compared pairwise against each other, with the agreement levels ranging between 96.3%–98.7%. The lowest level of mutual agreement is between the one-vs-rest and pairwise heuristics (96.3%), whereas the highest level is between the one-vs-rest and END heuristics (98.7%). Overall, all four heuristics agree with respect to the predicted lexeme in 3255 (95.6%) of the cases, so these results would also suggest that they all yield in the end relatively speaking very similar results.

Table 5.9. Pairwise comparisons of the lexemes predicted by each of the four polytomous logistic regression heuristics considered in this dissertation; absolute agreement figures supplemented with relative proportions in (parentheses).

THINK.multivariate.models_lexeme_selections.cross
THINK.multivariate.models_lexeme_selections.cross relative

Heuristics	pairwise	multinomial (simultaneous)	ensemble of nested dichotomies
one-vs-rest	3279 (96.3%)	3325 (97.7%)	3360 (98.7%)
pairwise	-	3313 (97.3%)	3312 (97.3%)
multinomial (simultaneous)	-	-	3344 (98.2%)

In general, the descriptive goodness-of-fit of the proper full model trained using the various heuristics with the entire data, measured in terms of relative decrease in deviance, can be considered relatively good since the associated measure R_L^2 ranges between 0.328–0.332.⁸⁸ Turning to using these fitted models to predict which lexeme should occur in a particular context, testing against the same data they were trained with, the different heuristics succeed at reaching a *Recall* rate of 63.81–63.89%, this relative difference corresponding in absolute terms in practice to no more than 10

⁸⁸ One should recall here that low R_L^2 values are overall the norm (as Hosmer and Lemeshow 2000: 167 note), and should not be compared as such with the corresponding R^2 statistic used in ordinary regression.

lexemes. Evaluated in terms of the reduction of prediction or classification error, the different heuristics succeed in beating the default choice of betting always for the most frequent lexeme *ajatella* (having a base-line success rate of almost every other time at $1492/3404=43.8\%$) with $\lambda_{prediction}=0.370-0.373$, and fare even better if measured against approaching on the long-run the overall relative proportions of all four selected THINK lexemes with $\tau_{classification}=0.490-0.493$. The ranges for the means of the corresponding validation figures are just slightly lower at $R_L^2=0.287-0.294$, $Recall=63.78-63.89\%$, $\lambda_{prediction}=0.355-0.357$ and $\tau_{classification}=0.478-0.480$, so the differences among the various heuristics are consistently small for all of the considered measures.

5.2.2 The lexeme-wise breakdown of the prediction results

We may next break down the predictions lexeme-wise, which allows us to estimate also their *Precision* in addition to *Recall*. As can be seen in Table 5.10, the lexemes certainly diverge from the lumped *Recall* level presented above. On the one hand, *ajatella* receives by far both the highest *Recall* (85.5%) and *Precision* (75.5%) values, which may again be attributed to its position as the most frequent of the selected THINK lexemes, accounting for close to one-half of the original occurrences. On the other hand, the three other, rarer lexemes fare less successfully, and while their prediction accuracy levels are broadly speaking quite similar with *Recall* ranging between 46.43–51.19% and *Precision* between 50.91–56.73%, it is interesting to note that the exact lexeme-wise values are almost in the same order as their original frequencies. Furthermore, these results all persist for the validated results produced with 1000-fold iterations using a simple bootstrap resampling on the entire data, presented in Table 5.11 further below.

Table 5.10. Lexeme-wise *Recall* and *Precision* in predicting outcomes in the entire original data ($n=3404$) using a single-fit proper full model (VI) with the one-vs-rest heuristic. THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected\$stats.lx

Lexeme	Original frequency	Relative original frequency (%)	Frequency of correct predictions	Recall (%)	Frequency of overall predictions	Precision (%)
ajatella	1492	43.8	1275	85.5	1758	72.5
miettiä	812	23.9	377	46.4	666	56.6
pohtia	713	21.0	365	51.2	624	58.5
harkita	387	11.4	182	47.0	356	51.1

Table 5.11. Lexeme-wise *Recall* and *Precision* in predicting outcomes in the entire original data ($n=3404$) using a proper full model (VI) fitted 1000-fold with the one-vs-rest heuristic and simple bootstrap resampling, values calculated assuming a normal distribution. THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se mantics_selected.1000\$stats.lx

Lexeme	Absolute original frequency	Mean frequency of correct predictions	Recall mean (%)	Absolute Std. Dev. of correct predictions	Recall Std. Dev. (%)	Mean frequency of overall predictions	Precision (%)
ajatella	1492	1274	85.4	17	1.13	1774	71.8
miettä	812	366	45.2	16	2.02	658	55.8
pohtia	713	353	49.5	18	2.56	621	56.9
harkita	387	178	46.0	10	2.47	351	50.8

We may further cross-tabulate the original lexemes against the predicted ones, presented for the single-fit model using one-vs-rest heuristic on the entire data in Table 5.12. As could be expected, for each original lexeme the most frequently predicted one is always the lexeme itself, with the corresponding proportion of such correct predictions equaling the lexeme-wise *Recall* values reported in Table 5.10 above. Likewise, for each predicted lexeme overall, the lexeme itself accounts for the largest proportion of original occurrences, in which case the corresponding (correct) proportion matches the lexeme-wise *Precision* values from Table 5.10 above. Focusing on the incorrect predictions instead and looking firstly from the original towards the predicted lexemes, *ajatella* would appear to be mistaken fairly rarely as any one of the other three THINK lexemes, in comparatively roughly equal proportions (3.9–4.2–6.5%). For its part, *miettä* has a relatively high chance of being predicted incorrectly as *ajatella* (28.9%), while to a lesser extent as *pohtia* (17.6%), and quite seldom as *harkita* (7.0%). In turn, *pohtia* has surprisingly close probabilities of being mistaken as either *ajatella* (20.3%) or *miettä* (20.2%), but it is more rarely confused with *harkita* (8.3%). Finally, *harkita* is quite often mixed up with *ajatella* (26.6%), but rather rarely though in roughly equal proportions with either *pohtia* (14.0%) or *miettä* (12.4%).

Consequently, all three rarer THINK lexemes are most often incorrectly predicted as *ajatella*, which may reflect some level bias in the setup of the overall polytomous model towards this most frequent one of the entire lot. As the constituent binary models in the one-vs-rest heuristic each contrast one of the lexemes against all the rest, there is simply much more negative evidence *against* the occurrence of each of the rarer lexemes individually (i.e., $n[-miettä]=2592$, $n[-pohtia]=2691$, and $n[-harkita]=3017$) than there is positive evidence *for* the occurrence of the most frequent lexeme ($n[ajatella]=1492$).

When switching next to the contrary perspective from the predicted towards the original lexemes, the highest proportion of incorrect predictions of *ajatella* is accounted originally for by *miettä* (13.4%), followed at some distance by *pohtia* (8.2%) and then *harkita* (5.9%). Next, among the mistaken predictions of *miettä*, *pohtia* (21.6%) ranks highest among the original, correct lexemes, with first *ajatella* (14.6%) and then *harkita* (7.2%) considerably lower down the line. In the case of incorrect predictions as *pohtia*, *miettä* (22.9%) would have been the correct lexeme most of the time, and to a clearly lesser but roughly equal extent either *ajatella* (9.9%) or *harkita* (8.7%). Lastly, approximately similar proportions of predictions of *harkita*

should rather be either *ajatella*, *mieltiä* or *pohtia* (16.0–16.3–16.6%). Thus, in comparison to the overall incorrect preference of *ajatella* demonstrated above, there is no single dominant lexeme which is mistaken for all the others. However, we may note that *mieltiä* and *pohtia* are mutually most often mistaken for each other, at an equal level of once in every five instances (21.6 vs. 22.9%).

Table 5.12. Distributions and proportions of predicted against original lexemes using a single-fit proper full model (VI) with the one-vs-rest heuristic on the entire data ($n=3404$); relative proportions of predicted lexemes out of original ones succeeded by the relative proportions of original lexemes among the predicted ones in parentheses as ($p_{\text{predicted}|\text{original}}|p_{\text{original}|\text{predicted}}$).
 THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se
 mantics_selected\$test.guess.mean
 THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se
 mantics_selected\$test.guess.rel

Original/ Predicted	ajatella	mieltiä	pohtia	harkita	$\Sigma(\text{orig.})$
ajatella	1275 (85.5 72.5%)	97 (6.5 14.6%)	62 (4.2 9.9%)	58 (3.9 16.3%)	1492
mieltiä	235 (28.9 13.4%)	377 (46.4 56.6%)	143 (17.6 22.9%)	57 (7.0 16.0%)	812
pohtia	145 (20.3 8.2%)	144 (20.2 21.6%)	365 (51.2 58.5%)	59 (8.3 16.6%)	713
harkita	103 (26.6 5.9%)	48 (12.4 7.2%)	54 (14.0 8.7%)	182 (47.0 51.1%)	387
$\Sigma(\text{predicted})$	1758	666	624	356	3404

Comparing the aforementioned values with those produced by the 1000-fold simple bootstrap of the same proper full model (VI) with the same one-vs-rest heuristic, the figures remain approximately the same, as is shown in Table 5.13 further below. At this stage, one could hypothesize that the lexeme-wise proportions of these mistaken predictions might be a proxy for semantic affinity and possibly even mutual interchangeability. This could be expected to hold especially if we take such similarity to be represented in the associated observable usage contexts, which is precisely what these predictions have been based on. Thus, the mutually highest confusion of *mieltiä* with *pohtia* could be taken as an indication of their close synonymy, supporting the conclusion to which I had come on the basis of manual scrutiny already in Arppe (2002).

Table 5.13. Mean distributions of predicted against original lexemes using a proper full model (VI) fitted 1000-fold with the one-vs-rest heuristic and simple bootstrap resampling on the entire data ($n=3404$); relative proportions of predicted lexemes out of original ones succeeded by the relative proportions of original lexemes among the predicted ones in parentheses as

$$(p_{\text{predicted}|\text{original}}|p_{\text{original}|\text{predicted}}).$$

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se
mantics_selected.1000\$test.guess.mean
THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se
mantics_selected.1000\$test.guess.rel

Original/ Predicted	ajatella	miettä	pohtia	harkita	$\Sigma(\text{orig.})$
ajatella	1274 (85.4 71.8%)	97 (6.5 14.7%)	63 (4.2 10.1%)	58 (3.9 16.5%)	1492
miettä	242 (29.7 13.6%)	367 (45.1 55.7%)	149 (18.4 24.0%)	55 (6.7 15.6%)	812
pohtia	154 (21.5 8.7%)	147 (20.6 22.3%)	353 (49.5 56.8%)	60 (8.4 17.1%)	713
harkita	105 (27.2 5.9%)	48 (12.3 7.2%)	56 (14.6 9.1%)	178 (46.0 50.8%)	387
$\Sigma(\text{predicted})$	1774	658	621	351	3404

5.2.3 Comparing the performance of models with different levels of complexity

Next, I will shift the focus from the performance of the various heuristics with respect to only one particular model, to the different types of models with varying levels of linguistic features and analytical complexity. In this, I will employ the one-vs-rest heuristic throughout on the basis of the arguments laid out earlier in Section 3.4.3. As can be seen in Table 5.14, increasing the number of feature categories and levels in linguistic analysis quite naturally has a positive impact on how much of the occurrences of the selected THINK lexemes can be accounted for. These results largely conform to those observed within the computational linguistic domain in, for example, classifying word senses on the basis of various combinations of different levels of automatic linguistic analysis (Lindén 2004). Starting at the simplest end, node-specific morphology (Model I), and somewhat surprisingly even if supplemented with verb-chain general morphological features (Model II), as well as extra-linguistic features alone (Model X), appear to have roughly equal (and low) explanatory power both in terms of fit with the original data as well as their added value in prediction. The *Recall* levels for these three models (I: 47.15%, II: 47.71%, and X: 47.21%) do not substantially rise above the base-line proportion of the most frequent THINK lexemes, *ajatella*, in the research corpus, being $1492/3404=43.8\%$. This is in fact reflected in the measures concerning the reduction of prediction error with $\lambda_{\text{prediction}}$ ranging 0.059-0.060-0.059, which indicate a minimal improvement in the results over always predicting the most frequent outcome class. In contrast, the measures for the reduction of classification error with these models are already clearly higher, with $\tau_{\text{classification}}$ ranging at 0.239-0.240-0.247, but among all the models considered here these values rank, nevertheless, as the lowest.

Syntactic argument types alone (Model III), without any of their semantic and structural subtypes, fare already slightly better. The fit with the original data is roughly equal to that achieved with the node-specific and verb-chain general

morphological features (Models I-II), and almost twice the corresponding value for extralinguistic features (Model X). As *Recall* with Model III increases to above the half-way-mark, the measures of prediction and classification error improve also accordingly, with $\lambda_{prediction}$ almost doubling in value in contrast to Models I-II and X; for $\tau_{classification}$ the absolute improvement is of a similar magnitude but lesser in relative terms. When morphological features concerning the entire verb-chain and the node-verb are combined with syntactic argument types (Model IV), the performance on the whole notches up noticeably. Now, the fit with the original data at $R_L^2=0.180$ is almost twice that of the morphological or syntactic arguments types alone (Models I-III), and over three times the level reached with extralinguistic features (Model X). Whereas *Recall* increases moderately to only 56.82%, especially the reduction of prediction error in comparison to syntactic argument types alone (Model III) roughly doubles, and also classification error reduces considerably, with $\lambda_{prediction}=0.231$ and $\tau_{classification}=0.378$.

If we further supplement the morphological and syntactic argument features with the semantic and structural classifications of the two most common and important arguments in the case of the studied THINK lexemes, namely, their AGENTS and PATIENTS (Model V), the results in terms of the descriptive fit of the model with the original data or prediction accuracy all improve again visibly. While *Recall* increases to 63.04%, the other measures grow less modestly by roughly one-third, as now $R_L^2=0.288$, $\lambda_{prediction}=0.342$, and $\tau_{classification}=0.468$. In contrast, adding further the subtypes for MANNER and TIME-POSITION arguments as well as the semantic classifications of verb-chains incorporated in the proper full model (VI) does not continue the improvement of the performance of the models at a rate similar to the immediately preceding additions in analytical intricacy and precision. Now, though descriptive fit has yet grown somewhat to $R_L^2=0.313$, on the predictive side *Recall* has increased by only one percent-unit to 64.6%, while the reduction of prediction error is modestly up with this model at $\lambda_{prediction}=0.370$ and $\tau_{classification}=0.490$.

It would appear that we are approaching some sort of upper limit, seemingly around a level of two-thirds accuracy in prediction, as to what can be achieved with the types of quite conventional linguistic analysis features applied in this study, concerning morphology, syntax and semantics within the immediate sentential context, since neither does the most complex model with the extended semantic classifications (Model VIII, with as many as 16 more semantic subtypes of syntactic arguments in comparison to Model VI) produce but quite minute improvements, with $R_L^2=0.325$, *Recall*=65.6%, $\lambda_{prediction}=0.388$, and $\tau_{classification}=0.504$. A similar conclusion was earlier noted in Arppe (2007) with a slightly differently selected extended variable set. Furthermore, dropping out the proper morphological verb-chain general features altogether but retaining the semantic classifications of verb-chains and combining these with the syntactic arguments as well as those among their semantic subtypes selected for the proper full model (VI), amounting to the feature set in model XI, results in a surprisingly small drop in performance, as $R_L^2=0.292$ with a *Recall*=63.1%, $\lambda_{prediction}=0.343$, and $\tau_{classification}=0.468$. Thus, the linguistic information coded in the morphological features, whether on the node-verb of the associated verb-chain in general, would appear to an essential extent be already incorporated in the syntactic and semantic argument structure.

As these results are clearly less than the performance levels achieved by Gries (2003b, *Recall*=88.9%, canonical *R*=0.821) and Bresnan et al. (2007, *Recall*=92%), even if achieved in simpler dichotomous settings, one possible avenue for improvement would be to add entirely new linguistic analysis categories such as longer-distance discourse factors as was done in these prior studies. However, the addition of a few extra-linguistic variables indicating medium and repetition in Arppe (2008) had no substantial effect (amounting then in practice only to an addition of 19 correctly classified selections). Likewise, the inclusion of the two extralinguistic features selected in this study, indicating the medium of usage (newspaper vs. Internet newsgroup discussion, and quoted fragments vs. body text), yield only small improvements of around one percent-unit in magnitude for the various performance measures, with $R_L^2=0.325$, *Recall*=65.57%, $\lambda_{prediction}=0.387$, and $\tau_{classification}=0.504$ for Model VII, and $R_L^2=0.337$, *Recall*=65.8%, $\lambda_{prediction}=0.391$, and $\tau_{classification}=0.507$ for Model IX. These results correspond in absolute terms to 33 more correctly classified lexeme selections in Model VII in comparison to Model VI, but only 7 in Model IX in comparison to Model VIII. Nevertheless, the results achieved by Inkpen and Hirst (2006: 25-27; see also Inkpen 2004: 111-112), with over 90 percent accuracy in correctly selecting a synonym from several multiple-lexeme sets, would suggest that the choices can in fact be highly exactly modeled. However, this required explanatory variables indicating “nuances” such as denotational microdistinctions as well as expressive ones concerning the speaker’s intention to convey some attitude, in addition to the sought-after style, which are not necessarily explicitly evident in the immediate sentential context nor easily amenable to accurate automated extraction (Edmonds and Hirst 2002: 128, cf. Hanks 1996: 90, 97).

The current performance plateau may result from technical restrictions related to the application of the one-vs-rest heuristic in particular, and on the basis of the similarities in the performance of all the heuristics demonstrated above, of polytomous logistic regression in general, to the more complex, multiple-outcome setting in this study. This may also result to some extent from the exclusion of interaction terms among the explanatory variables included in all the Models I-XI presented above, due to restrictions set by the size of the available data. But this might also reflect genuine synonymy, or at least some extent of interchangeability in at least some contexts, which the current analysis variables cannot (possibly ever) get an exact hold of (cf. Gries 2003b: 13-16). Even more radically we may interpret such (varying degrees of) interchangeability as evidence rather for inherent variability in language, following Bresnan (2007).

Thus, though the individual linguistic choices, associated with some contexts represented as linguistic explanatory variables, have to be discrete for each instance of usage by an individual person at a particular place and time, over longer stretches of language usage such outcomes as studied here may turn out to be probabilistic instead. That is, the workings of a linguistic system, represented by the range of variables according to some theory such as the ones used in this dissertation, and its resultant usage need not in practice be categorical, following from exception-less rules, but may exhibit degrees of potential variation which becomes evident over repeated use, manifested in, for example, corpora the likes of those used here.

Table 5.14. The descriptive and predictive properties of the various types of Models (I-XI) with different compositions of explanatory variables, based on the single-fit training and testing of each model with the one-vs-rest heuristic data on the entire data ($n=3404$).

Model index	Feature set composition	Recall (%)	R_L^2	$\lambda_{prediction}$	$\tau_{classification}$
I	Only node-specific morphological features (26)	47.15	0.094	0.059	0.239
II	Verb-chain general morphological features (10) as well as those node-specific features which are not subsumed by the verb-chain general ones (17)	47.71	0.100	0.069	0.247
III	Syntactic argument types, <i>without</i> their semantic and structural classifications	50.18	0.098	0.113	0.282
IV	Verb-chain general morphological features (10) and non-subsumed node-specific morphological features (17) together with syntactic argument types (17), the latter again without their subtypes	56.82	0.180	0.231	0.378
V	Verb-chain general features (10), the most common semantic classifications of AGENTS and PATIENTS with the less frequent subtypes collapsed together (12), and the other syntactic argument types alone (15)	63.04	0.288	0.342	0.468
VI	Proper full model with verb-chain general morphological features (10) and their semantic classifications (6) together with syntactic argument types alone (10) or their selected or collapsed subtypes (20)	64.60	0.313	0.370	0.490
VII	Proper full model with verb-chain general morphological features (10) and their semantic classifications (6) together with syntactic argument types alone (10) or their subtypes (20) as well as extra-linguistic features (2)	65.57	0.325	0.387	0.504
VIII	Extended full model with verb-chain general morphological features (10) and their semantic classifications (9) together with syntactic argument types (5) and all their subtypes exceeding the minimum frequency threshold (38)	65.60	0.325	0.388	0.504
IX	Extended full model with verb-chain general morphological features (10) and their semantic classifications (9) together with syntactic argument types (5) and all their subtypes exceeding the minimum frequency threshold (38) as well as extra-linguistic features (2)	65.80	0.337	0.391	0.507
X	Extralinguistic features alone (2)	47.21	0.057	0.060	0.240
XI	Syntactic argument types alone (10) or their selected or collapsed subtypes (20), together with semantic classifications of verb chains (6) but <i>without</i> any morphological features	63.10	0.292	0.343	0.468

The diminishing effect of increasing analytical intricacy and the number of explanatory variables noted above, in addition to earlier observations of the apparent indifference to the exact composition of the variable set regarding the results (Arppe 2007), raise the somewhat heretical question of whether comparable results could be achieved by randomly selecting the set of variables incorporated in a model. Consequently, I decided to try out and observe what would happen if a variable set equaling in size (46) that included in the proper full model (VI) would be randomly sampled from the variables (62 in all) incorporated in the substantially larger extended full model (VIII). The somewhat surprising results of repeating this process 100 times are presented in Table 5.15. Though the average results are positioned between Models VI and V, the best randomly selected variable set (listed in Table R.10 in Appendix R) performs almost as well as the proper full model (VI), with $R_L^2=0.303$, $Recall=65.4\%$, $\lambda_{prediction}=0.367$, and $\tau_{classification}=0.488$, and even the worst such random variable set comes pretty close to Model V, with $R_L^2=0.203$, $Recall=65.8\%$, $\lambda_{prediction}=0.192$, and $\tau_{classification}=0.346$.

While the best such random model had as many as 31 (67.4%) variables in common with the proper full model (VI), the worst random model was just three variables worse off at 28 correspondingly common variables (60.9%). In between themselves, the best and worst performing random variable sets had mutually 32 variables in common (69.6%). Thus, it would seem that the entire considered variable set is interrelated in manifold ways and the different features rather represent different facets of the studied phenomenon than are fully distinct from each other. This would also suggest that some more abstract, as of yet unidentified variables, which may possibly not be manifested in singular words in the context but could rather concern the entire argument structure, might lie behind the more explicit ones now under consideration.⁸⁹ Such interrelationships and the posited underlying, more profound variables could be studied and identified statistically by using, for example, cluster or principal components analysis, as was noted earlier in Section 3.4.2, aggregating the current variable set into a smaller but more abstract one.

Table 5.15. The mean, maximum, and minimum descriptive and predictive properties of 100 models, for which each 46 explanatory variables were selected randomly from the 62 features in the extended full model (VIII), based on the single-fit training and testing of each model with the one-vs-rest heuristic data on the entire data ($n=3404$).

THINK.multivariate.one_vs_rest.46_random_variables.100\$variable.results.specific

Feature set composition	Recall (%)	R_L^2	$\lambda_{prediction}$	$\tau_{classification}$
Random variable sets on average	60.75 (2.24)	0.261 (0.025)	0.301 (0.040)	0.435 (0.032)
“Best” random variable set (31 features in common with proper full model)	64.42	0.303	0.367	0.488
“Worst” random variable set (28 features in common with proper full model)	54.61	0.203	0.192	0.346
Quartiles (25%/75%)	59.14/63.04	0.239/0.278	0.272/0.342	0.411/0.468

⁸⁹ This possible interpretation has been suggested to me by Professor Lauri Carlson.

5.3 Relative lexeme-wise weights of feature categories and individual features

5.3.1 Overview of the lexeme-specific feature-wise odds

We can now shift the focus to the individual explanatory variables and start off by evaluating for the final (proper) full model (VI) what are the average weights of the various variable categories. As can be seen in Table 5.16, when considering only significant odds either way, syntactic arguments coupled with a semantic classification are clearly the most distinctive group of features with a mean aggregate odds (based on the *absolute* values of the underlying log-odds⁹⁰) of 4.08, whether in favor (3.64) or against (0.22~1:4.55) the occurrence of a lexeme. The semantic characterizations of the verb chain have the second-highest impact (3.09), followed relatively closely by syntactic arguments alone (3.06), without any semantic or structural subtypes as is the case with the rarer arguments. Morphological features pertaining to both the node-verb and the possibly associated verb-chain have the least overall impact (2.12). Counting in all estimated odds, including the nonsignificant ones, the ranking order of the feature categories remains the same, though the differences between them in terms of their mean odds become greater.

Furthermore, it appears that mean (significant) odds *against* the occurrence of a lexeme are for each feature category stronger than those in favor, being 1:2.22 vs. 2.02:1 for verb-chain morphology, 1:4.17 vs. 2.60:1 for verb-chain semantics, 1:3.45 vs. 2.67:1 for syntactic argument types alone, and 1:4.55 vs. 3.64:1 for combinations of syntactic arguments with their semantic subtypes. This may follow from the empirical fact that especially for the rarer THINK lexemes their chances of occurrence are considerably fewer than their nonoccurrence (including the combined occurrences of all three other lexemes), and thus the “odds” are in general more against the occurrence of these lexemes than in their favor.

Table 5.16. Average weights of the different categories of explanatory feature variables, calculated firstly for significant odds in the overall polytomous regression model (VI) attained with the one-vs-rest heuristic on the entire data, and secondly for all estimated odds including nonsignificant ones (in parentheses).

Feature variable category	Mean odds in favor	Mean odds against	Mean aggregate odds
Verb chain morphology	2.02 (1.52)	0.45~1:2.22 (0.66~1:1.52)	2.12 (1.52)
Verb chain semantics	2.60 (1.73)	0.24~1:4.17 (0.16~1:6.25)	3.09 (3.45)
Syntactic argument types (alone)	2.67 (1.87)	0.29~1:3.45 (0.47~1:2.13)	3.06 (2.01)
Syntax arguments + semantic/structural subtypes	3.64 (2.68)	0.22~1:4.55 (0.06~1:17)	4.08 (7.93)

⁹⁰ In the calculation of the overall aggregate odds, in the case of odds *against* a lexeme (all which are $e^{\beta[L|exeme]|F[feature]} < 1$), e.g., 0.5, I have used their inverse value, i.e., $1/0.5=2$ in this particular case. When one uses the underlying log-odds values, which for any odds *against* a lexeme would be negative ($\beta[L|F] < 0$), the aforementioned procedure corresponds to taking the absolute values of the log-odds. The mean aggregate odds are then attained by calculating first the mean of the absolute log-odds, i.e., $=\bar{x}[\beta(L|F)]$, followed by raising e to the value of the attained mean, i.e., $e^{\bar{x}[\beta(L|F)]}$.

The individual lexeme-wise odds for all the feature variables incorporated in the proper full model (VI) are presented in its entirety in Table 5.17 below. These results can now be scrutinized from two perspectives, either lexeme-wise or feature-wise. From the lexeme-wise angle, we may rank the features per each lexeme as to how much they either increase the odds ($e^{\beta[L|F]} > 1 \sim \beta[L|F] > 0$) or decrease the odds ($e^{\beta[L|F]} < 1 \sim \beta[L|F] < 0$) of the particular lexeme occurring. At the same time, we can also note which features are neutral with respect to each lexeme, that is, features for which their lexeme-specific odds do not statistically significantly diverge from 1.0. The number of significant odds per lexeme appears to be associated with the overall frequency of the lexeme, as for the most frequent *ajatella* 32 features overall exhibit significant odds either in favor of or against it⁹¹, while the respective figures for the rarer lexemes are 22 for *miettiä*, 20 for *pohtia*, and 13 for *harkita*. More specifically, among the significant odds for each lexeme, 15 are in favor of and 17 against the occurrence of *ajatella*, whereas the corresponding figures are 14 vs. 8 for *miettiä*, 12 vs. 8 for *pohtia*, and 6 vs. 7 for *harkita*. Thus, the balance of features in favor of or against a lexeme varies, with *miettiä* and *pohtia* having more features in their favor, while *ajatella* and *harkita* have more features against their occurrence, relatively speaking. Furthermore, the number of neutral, nonsignificant features per each lexeme also varies, being 15 for *ajatella*, 25 for *miettiä*, 27 for *pohtia*, and 34 for *harkita*.

⁹¹ This corresponds with Divjak and Gries' (2006: 43) result that the most frequent of the near-synonymous Russian TRY verbs they studied, *probovat'*, also had the largest number of different contextual features (i.e., *ID tags* in their parlance) occurring with it (at least once).

Table 5.17. Odds of the proper full polytomous logistic regression model (VI) fitted using the one-vs-rest heuristic, with each of the studied THINK lexemes pitted against the others at a time; odds against any lexeme, i.e., $e^{\beta(L|F)} < 1$, supplemented by the corresponding ratio, i.e., $1:1/e^{\beta(L|F)} = 1:e^{-\beta(L|F)}$, e.g., 0.5~1:2; significant lexeme-wise odds in **boldface**; nonsignificant odds in (parentheses); features with at least one lexeme with significant odds in **boldface**.

Feature/Lexeme	ajatella	miettä	pohtia	harkita
SX AGE.SEM_GROUP	0.2~1:5	0.52~1:1.9	4.2	(1.1)
SX AGE.SEM_INDIVIDUAL	(0.85~1:1.2)	(0.98~1:1)	(1.6)	(0.69~1:1.5)
SX_CND	0.46~1:2.2	(1.2)	(0.57~1:1.7)	2.9
SX CV	0.48~1:2.1	2.3	(0.84~1:1.2)	(0.81~1:1.2)
SX DUR	0.12~1:8.4	3.4	(1.3)	(1)
SX FRQ	0.38~1:2.6	1.7	(0.79~1:1.3)	(1.7)
SX GOA	3.8	(0.56~1:1.8)	(0.57~1:1.8)	0.21~1:4.7
SX LOC	0.26~1:3.9	(0.93~1:1.1)	3.7	0.46~1:2.2
SX_LX_että_CS.SX PAT	2.6	0.52~1:1.9	0.5~1:2	0.25~1:4
SX MAN.SEM AGREEMENT	16	0.07~1:14	0.22~1:4.5	(0~1:7e ⁶)
SX MAN.SEM FRAME	2.4	0.28~1:3.6	(1.3)	0.27~1:3.8
SX MAN.SEM GENERIC	23	0.15~1:6.8	(0~1:5e ⁶)	(0~1:9e ⁶)
SX MAN.SEM JOINT	0.37~1:2.7	2.1	(0.78~1:1.3)	(1.5)
SX MAN.SEM NEGATIVE	4	(0.56~1:1.8)	0.22~1:4.6	(0.58~1:1.7)
SX MAN.SEM POSITIVE	(0.71~1:1.4)	(0.99~1:1)	(0.82~1:1.2)	1.8
SX META	(0.83~1:1.2)	(1)	(0.8~1:1.2)	1.6
SX PAT.DIRECT QUOTE	0.013~1:75	3	8.1	(0~1:8.1e ⁶)
SX PAT.INDIRECT QUESTION	0.07~1:14	4.2	2.8	(0.82~1:1.2)
SX PAT.INFINITIVE	5.3	(0~1:4e ⁶)	(0.21~1:4.7)	(1.4)
SX PAT.PARTICIPLE	5.3	(0~1:4e ⁶)	(0.3~1:3.3)	(1.1)
SX PAT.SEM ABSTRACTION	0.25~1:4.1	1.5	4.1	(1)
SX PAT.SEM ACTIVITY	0.14~1:7.1	(0.77~1:1.3)	1.6	9
SX PAT.SEM COMM...	0.1~1:9.6	2.8	3	(1.8)
SX PAT.SEM_EVENT	(1.4)	(0.97~1:1)	(0.98~1:1)	(0.34~1:3)
SX PAT.SEM_INDIV..._GROUP	2.7	0.52~1:1.9	0.3~1:3.4	(0.87~1:1.2)
SX_QUA	(0.69~1:1.5)	2.6	(0.75~1:1.3)	0.33~1:3
SX_RSN_PUR	0.43~1:2.3	(1.1)	(1.3)	(1.6)
SX_SOU	3.1	(0.76~1:1.3)	0.29~1:3.5	0.13~1:7.5
SX_TMP.SEM DEFINITE	0.4~1:2.5	(0.97~1:1)	2.3	(0.76~1:1.3)
SX_TMP.SEM INDEFINITE	0.57~1:1.7	1.5	(0.97~1:1)	(1.2)
SX_VCH.SEM ACCIDENTAL	5.6	(0.44~1:2.3)	(0.48~1:2.1)	(0~1:1e ⁷)
SX_VCH.SEM EXTERNAL	2.5	(0.8~1:1.3)	(0.73~1:1.4)	(0.91~1:1.1)
SX_VCH.SEM NECESSITY	0.35~1:2.9	2	(0.96~1:1)	(1.4)
SX_VCH.SEM_POSSIBILITY	(1.2)	(1.1)	(0.82~1:1.2)	(1.2)
SX_VCH.SEM TEMPORAL	0.26~1:3.8	1.8	2.4	0.15~1:6.5
SX_VCH.SEM_VOLITION	(0.64~1:1.6)	(1.6)	(1)	(0.64~1:1.6)
Z_ANL_COVERT	(1.1)	(1.2)	(0.77~1:1.3)	(0.79~1:1.3)
Z_ANL_FIRST	(0.86~1:1.2)	(1.8)	0.29~1:3.5	(1.9)
Z_ANL_IND	2	(0.67~1:1.5)	(0.81~1:1.2)	(0.81~1:1.2)
Z_ANL_KOND	(1.3)	0.54~1:1.9	(0.7~1:1.4)	2.3
Z_ANL_NEG	2.1	(0.72~1:1.4)	0.48~1:2.1	(1.1)
Z_ANL_PASS	(0.63~1:1.6)	(0.89~1:1.1)	1.9	(1.1)
Z_ANL_PLUR	(1.1)	0.59~1:1.7	1.6	(1.2)
Z_ANL_SECOND	(0.69~1:1.5)	2.4	0.42~1:2.4	(0.68~1:1.5)
Z_ANL_THIRD	(0.63~1:1.6)	(1.3)	(0.99~1:1)	(1.6)
Z_PHR_CLAUSE	(1.1)	(0.59~1:1.7)	(0.87~1:1.1)	(2)

5.3.2 Lexeme-wise analysis of the estimated odds

In general, whereas features with odds in favor of the occurrence of a lexeme can be expected to genuinely occur in conjunction with the particular lexeme, features with odds against the occurrence of a lexeme rather position the lexeme in contrast with the entire lexical set studied at the same time, and can thus be expected instead to be found predominantly in the contexts of one or more of the other scrutinized lexemes. Looking at the individual lexemes, we may see in Table 5.18 presenting the ordering of the features for *ajatella* with respect to their odds that particularly the GENERIC and AGREEMENT types of MANNER increase the odds of *ajatella* occurring substantially, at ratios of 23:1 and 16:1, respectively, followed at quite a distance by the ACCIDENTAL verb chain construction (5.6:1), both INFINITIVES and PARTICIPLES as PATIENTS (5.3:1), NEGATIVE evaluations of MANNER (4:1), GOAL (3.8:1) and SOURCE (3.1:1) arguments, INDIVIDUALS and GROUPS combined (2.7:1) as well as *että*-clauses (2.6:1) as PATIENTS, an indication of EXTERNAL cause in the verb chain, a FRAME as MANNER (2.4:1), and finally NEGATION (2.1:1) or the INDICATIVE mood (2:1) morphologically manifested in the verb-chain.

In contrast, either a DIRECT QUOTE or an INDIRECT QUESTION as a PATIENT in the context tip the scales considerably against the occurrence of *ajatella*, at ratios of 1:75 and 1:14, respectively, as is the case to a lesser extent also with expressions or media of COMMUNICATION (1:9.6) as PATIENT, DURATION as an argument (1:8.4), ACTIVITIES as PATIENT (7.1:1), GROUPS as AGENT (1:5.0), ABSTRACTIONS as PATIENT (4.1:1), LOCATION arguments (1:3.8), an expression of TEMPORALITY (1:2.9) or NECESSITY (1:3.1) in the verb chain, the JOINT subtype of MANNER (1:2.7), a FREQUENCY argument (1:2.6), a DEFINITE expression of TIME-POSITION (1:2.5), REASON or PURPOSE combined (1:2.3) or CONDITION (1:2.2) as an argument or a CO-ORDINATED VERB (1:2.0), and lastly, also an INDEFINITE expression of TIME-POSITION (1:1.7).

Table 5.18. Features with significant odds either in favor of or against *ajatella*.

Odds in favor (15)	Odds against (17)
SX_MAN.SEM_GENERIC (23)	SX_PAT.DIRECT_QUOTE (0.013~1:75)
SX_MAN.SEM_AGREEMENT (16)	SX_PAT.INDIRECT_QUESTION (0.07~1:14)
SX_VCH.SEM_ACCIDENTAL (5.6)	SX_PAT.SEM_COMMUNICATION (0.1~1:9.6)
SX_PAT.INFINITIVE (5.3)	SX_DUR (0.12~1:8.4)
SX_PAT.PARTICIPLE (5.3)	SX_PAT.SEM_ACTIVITY (0.14~1:7.1)
SX_MAN.SEM_NEGATIVE (4)	SX_AGE.SEM_GROUP (0.2~1:5)
SX_GOA (3.8)	SX_PAT.SEM_ABSTRACTION (0.25~1:4.1)
SX_SOU (3.1)	SX_LOC (0.26~1:3.9)
SX_PAT.SEM_INDIV..._GROUP (2.7)	SX_VCH.SEM_TEMPORAL (0.26~1:3.8)
SX_LX_että_CS.SX_PAT (2.6)	SX_VCH.SEM_NECESSITY (0.35~1:2.9)
SX_VCH.SEM_EXTERNAL (2.5)	SX_MAN.SEM_JOINT (0.37~1:2.7)
SX_MAN.SEM_FRAME (2.4)	SX_FRQ (0.38~1:2.6)
Z_ANL_NEG (2.1)	SX_TMP.SEM_DEFINITE (0.4~1:2.5)
Z_ANL_IND (2)	SX_RSN_PUR (0.43~1:2.3)
	SX_CND (0.46~1:2.2)
	SX_CV (0.48~1:2.1)
	SX_TMP.SEM_INDEFINITE (0.57~1:1.7)

Similar assessments can be done for each of the THINK lexemes included in the analysis, and are presented in full in Tables R.12-15 in Appendix R. Nevertheless,

features at least doubling the odds either in favor of or against the occurrence of each of the three other THINK lexemes are also mentioned here, i.e., with either $odds \geq 2$ or $odds \leq 0.5$. For *miettiä*, the strongest odds in its favor are in conjunction with an INDIRECT QUESTION as a PATIENT (4.2:1), followed by DURATION as a syntactic argument (3.4:1), a DIRECT QUOTE (3.0:1) or an expression or medium of COMMUNICATION (2.8:1) as PATIENT, QUANTITY (2.6:1) as an argument, SECOND person expressed morphologically in the verb-chain (with the odds 2.4:1, co-occurring roughly half of time the IMPERATIVE mood), a CO-ORDINATED VERB (2.3:1), the JOINT subtype of MANNER (2.1:1), and finally the expression of NECESSITY in the verb-chain (2:1). In contrast, the strongest odds against *miettiä* are the AGREEMENT (0.07~1:14), GENERIC (0.15~1:6.8), and FRAME (0.28~1:3.6) subtypes of MANNER.

For *pohtia*, the strongest odds in its favor are in conjunction with a DIRECT QUOTE as PATIENT (8.1:1), followed by a GROUP as an AGENT (4.2:1), and ABSTRACTION as a PATIENT (4.1:1), LOCATION (3.7:1) as an argument, expressions or media of COMMUNICATION (3.0:1) or an INDIRECT QUESTION (2.8:1) as a PATIENT, a TEMPORAL expression in the verb-chain (2.4:1) or a DEFINITE expression of TIME-POSITION (2.3:1). To the contrary, AGREEMENT or (0.22~1:4.5) a NEGATIVE evaluation (0.22~1:4.6) as MANNER, SOURCE (0.29~1:3.5) as an argument, FIRST person (0.29~1:3.5) expressed morphologically in the verb-chain, either a human INDIVIDUAL or GROUP as PATIENT (0.3~1:3.4), SECOND person (0.42~1:2.4) or NEGATION (0.48~1:2.1) as well as an *että*-clause ‘that’ as PATIENT (0.5~1:2) exhibit the strongest odds against *pohtia*.

Finally, with respect to *harkita*, an ACTIVITY as PATIENT (9:1), CONDITION (2.9:1) in general as an argument, and CONDITIONAL mood in the verb-chain have the strongest odds in favor of this lexeme. As the strongest odds against *harkita* are SOURCE (0.13~1:7.5) as an argument, TEMPORALITY (0.15~1:6.5) expressed in the verb-chain, GOAL (0.21~1:4.7) as an argument, an *että*-clause as a PATIENT (0.25~1:4), FRAME as MANNER (0.27~1:3.8), as well as QUANTITY (0.33~1:3) or LOCATION (0.46~1:2.2) in general as syntactic arguments.

5.3.3 Feature-wise analysis of the estimated odds

In contrast, from the feature-wise viewpoint we may be interested in which of the lexemes have the strongest odds in favor of (>1) or against (<1) occurring in conjunction with each selected individual feature or groups of related features, and for which lexeme(s) the odds are neutral. In all, there are two features for which *all* the lexeme-wise odds are significant, namely, an *että* ‘that’ clause as a PATIENT, preferring *ajatella* and dispreferring all the rest, and TEMPORALITY expressed in the verb chain, preferring both *miettiä* and *pohtia* while dispreferring *ajatella* and *harkita*. In contrast, there were 7 features for which all lexeme-wise odds are nonsignificant, namely, INDIVIDUALS as AGENT, EVENTS as PATIENT, POSSIBILITY and VOLITION expressed in the verb chain, the COVERTNESS of the AGENT, THIRD PERSON expressed morphologically in the verb-chain, and usage as a CLAUSE-EQUIVALENT form. In between these extreme ends, there were 10 features for which only one of the four THINK lexemes had significant lexeme-specific odds, while 16 features had a significant effect on exactly two lexemes, and 11 features on precisely three lexemes.

Focusing specifically on the feature-wise results for AGENT-related syntactic, semantic and verb-chain general morphological features presented in Table 5.19, we can firstly see that a human INDIVIDUAL as AGENT is neutral with respect to all studied THINK lexemes, whereas a human GROUP in the same argument slot is considerably more discriminatory, in that it has significant odds in favor of *pohtia* and against both *miettiinä* and *ajatella*, while remaining neutral with respect to *harkita*. Among the three persons, the FIRST PERSON exhibits only a significant dispreference for occurring with *pohtia*, whereas it is neutral for the three other THINK lexemes, leaving no lexeme with significant odds in their favor. The SECOND PERSON is more selective as it has significant odds in favor of *miettiinä* and against *pohtia*, staying relatively neutral in the case of both *ajatella* and *harkita*. However, the THIRD PERSON shows no significant odds either in favor of or against any of the four lexemes.

In terms of morphological (verbal) number, the PLURAL has significant odds for *pohtia* and against *miettiinä*, with *harkita* and *ajatella* as neutral, which is somewhat similar to the preferences of GROUP AGENTS and a mirror image of the SECOND PERSON.⁹² The impersonal PASSIVE voice also exhibits significant odds for *pohtia*, but it remains only neutral in conjunction with the three other lexemes. In this respect it is somewhat unexpected that not having an OVERT agent (denoted by the tag Z_ANL_COVERT) is neutral with respect to all four THINK lexemes. Finally, a LOCATION as an argument has significant odds in favor of *pohtia* and against *ajatella* and *harkita*, while it stays neutral for *pohtia*.

Altogether, I interpret these results to entail that human GROUPS, whether explicitly indicated as collectives or countable groups of individuals, indirectly referred to via a LOCATION, or unidentified and impersonal as is implied with PASSIVE voice, were attracted to *pohtia*, whereas *miettiinä* and *ajatella* appear repulsed by these same characteristics. In line with this, discourse-proximal reference to either individual speaker(s) or addressee(s) in the FIRST and SECOND PERSONS, respectively, would shirk *pohtia*, and furthermore in the case of the SECOND PERSON instead be specifically associated with *miettiinä*. As *pohtia* has overall clearly the largest proportion of significant odds in favor of or against this subset of AGENT-related features (6), in comparison to the other three THINK lexemes (3 for *miettiinä*, 2 for *ajatella*, and only one for *harkita*), it would seem the most specialized and distinguished one among this lexeme set with respect to the type of agency it represents, which would also conform with the general *post hoc* hypotheses proposed earlier in Section 4.1.2.

⁹² Comparing these results against those for the six specific features combining both person and number which were used in the example in Section 3.4.5, contrasting *ajatella* against the rest, we may note that the FIRST PERSON SINGULAR, FIRST PERSON PLURAL as well as THIRD PERSON PLURAL, all with significant odds in favor of *ajatella*, have not persisted as significant among the four more generalized person and number features applied here.

Table 5.19. The feature-wise sorting of the studied THINK lexemes per the two semantic subtypes of AGENT as well as the related verb-chain general morphological features and the superficially unrelated LOCATION argument into ones with significant odds in favor of, neutral (nonsignificant), and significant odds against the occurrence of each lexeme.

Contextual feature	Lexemes with significant odds in favor	Lexemes with neutral odds	Lexemes with significant odds against
SX_AGE.SEM_INDIVIDUAL	-	pohtia (1.6), miettä (0.98), ajatella (0.85), harkita (0.69)	-
SX_AGE.SEM_GROUP	pohtia (4.2)	harkita (1.1)	miettä (0.52), ajatella (0.2)
Z_ANL_FIRST	-	harkita (1.9), miettä (1.8), ajatella (0.86)	pohtia (0.29)
Z_ANL_SECOND	miettä (2.4)	ajatella (0.69), harkita (0.68)	pohtia (0.42)
Z_ANL_THIRD	-	harkita (1.6), miettä (1.3), pohtia (0.99), ajatella (0.63)	-
Z_ANL_PLUR	pohtia (1.6)	harkita (1.2), ajatella (1.1)	miettä (0.59)
Z_ANL_PASS	pohtia (1.9)	harkita (1.1), miettä (0.89), ajatella (0.63)	-
Z_ANL_COVERT	-	miettä (1.2), ajatella (1.1), harkita (0.79), pohtia (0.77)	-
SX_LOC	pohtia (3.7)	miettä (0.93)	harkita (0.46), ajatella (0.26)

We can now compare these results with an earlier multimethodological study (Arppe and Järvikivi 2007b) which combined both corpus and experimental data concerning the AGENT types and the associated person/number features, and which focused only on the lexeme pair *miettä* and *pohtia*. Within the more complex syntactic-semantic network and the larger group of THINK lexemes considered in this study, it is interesting to note that the contrasts observed between *miettä* and *pohtia* shift somewhat, but are nonetheless essentially upheld. As concluded in the combined results in the earlier study, a human GROUP as an AGENT has strong and significant odds in favor of *pohtia* and against *miettä*, the latter which was in particular evident in the acceptability rating experiments of the former study. With respect to human INDIVIDUALS as AGENT, the results in this study conform to the overall conclusion in the prior study that there is no significant difference between the two lexemes for this feature combination.

Furthermore, whereas the corpus-based results in the prior study indicated a strong positive association between FIRST PERSON SINGULAR and *miettä*, and a negative one with *pohtia*, in this study the result stays the same for *pohtia*, while the effect with respect to *miettä* has become neutral. It might be conceivable, however, that the

association between FIRST PERSON and *miettiinä* still remains but has simply been surpassed by an even stronger preference for the same lexeme by the SECOND PERSON. Nevertheless, it would be interesting to find out whether this dispreference for *pohtia* with respect to the FIRST PERSON would diminish also in an acceptability rating experiment covering all the four THINK lexemes and all person/number features, similar to what was observed in such an experiment in the earlier study.

5.3.4 Comparison with the univariate results

More generally, we can also compare these multivariate results with the univariate ones presented in Section 4.1 and Appendix N. Firstly, at the feature level, one might correctly assume that higher overall levels of association between the selected THINK lexemes and each particular feature would correlate to some degree, at least in terms of ranking order, with the overall strength of lexeme-wise odds per each feature. Indeed, the Spearman rank order correlation coefficient between the lexeme-wise U_{FL} association measures and the mean aggregated odds per each feature calculated over the four lexemes (presented in full in Table R.16 in Appendix R) is very high, at $r_s=0.827$, which would suggest that the lexeme-wise association values acquired in the univariate analysis are a relatively good indicator of which features will turn out to be significant in the multivariate analysis, and of how strong their relative importance will be relative to the other features. Secondly, one could further very well entertain the idea that the strength of lexeme-wise deviation from a homogeneous distribution assessed for all features in the univariate analyses with standardized Pearson residuals would correlate at least to some extent with the lexeme-wise odds assigned for these same features in the multivariate logistic regression results presented here. However, it appears that there is practically no correlation, at least as to the strength of these values, since the ranked Spearman coefficient is overall $r_s=-0.045$ for all feature-lexeme pairings regardless of their significance in the multivariate results, and $r_s=-0.073$ if considering only the features with significant odds. Neither does applying the same evaluation per each lexeme produce evidence of stronger relationships, as $r_s(\text{ajatella})=-0.166$, $r_s(\text{miettiinä})=0.092$, $r_s(\text{pohtia})=-0.134$, and $r_s(\text{harkita})=-0.065$ (without excluding feature-lexeme pairings with nonsignificant odds).

If we simply look at the directions of the preferences indicated in either the univariate or the multivariate analyses, summarized in Table 5.20, we may first note that there are no reversals, that is, cases in which a positive or negative preference in the univariate results would receive in the multivariate analyses odds in the opposite direction. In contrast, 40 of the positive lexeme-wise preferences (+) are also assigned significant odds in favor of (>1), and 36 of the dispreferences (–) are assigned odds against (<1) the occurrence of the particular lexeme, while 55 instances of lexemes neutral (0) with respect to a given feature remain so also in the multivariate analysis (i.e., $odds \approx 1$). In sum, this means that a clear majority of 131 (71.2%) feature-lexeme associations retain their directions of preference (or lack thereof). Overall, with respect to the direction of lexeme-wise preferences and odds there is a strong association, as the distribution in Table 5.20 is as a whole firstly significant with $P(df=4)=6.07e^{-32}$. Furthermore, using the asymmetric Uncertainty Coefficient we may note that the multivariate directions for lexeme-wise preferences are a somewhat better predictor of the corresponding directions of univariate preferences, with

$U_{Univariate|Multivariate}=0.426$ and $U_{Multivariate|Univariate}=0.389$, but in any case these association values for both directions can be characterized as quite strong.

Table 5.20. Comparison of the lexeme-specific preferences from the univariate analysis with the lexeme-specific odds from the multivariate analysis; ‘+’ denoting a positive preference and (>1) significant odds in favor of a lexeme in conjunction with a feature, ‘-’ a dispreference and (<1) significant odds against a lexeme in such a context, and ‘0’ a neutral relationship and (≈ 0) nonsignificant odds in the two respective analyses.

Univariate/Multivariate	>1	≈ 0	<1	Σ
+	40	3	0	43
0	16	55	30	101
-	0	4	36	40
Σ	46	62	66	184

The strongest changes between the two levels of analyses concern cases in which neutral features turn out to have significant odds for a given lexeme, or significant preferences or dispreferences which lose their relative importance when their impact is considered comparatively together with all the other selected features. The former set of features with a shift from neutral univariate association to significant odds against a lexeme includes GOAL as an argument and FRAME as MANNER in conjunction with *harkita*, as well as the CONDITIONAL mood and PLURAL number with *miettiä*, whereas previously neutral features which end up having instead significant odds in favor of a lexeme covers expressions and media of COMMUNICATION as PATIENT with *miettiä*, an indication of EXTERNAL cause in a verb-chain with *ajatella*, and PLURAL number together with *pohtia*. The complete underlying univariate and multivariate preference patterns as well as the standardized Pearson residuals and odds on which they are based are presented in Tables R.17 and R.18 in Appendix R.

5.3.5 Comparison with descriptions in current dictionaries

Next, we can once more compare the corpus-based evidence, now in light of the multivariate analysis, with respect to the exposition of various features among the example sentences for the four THINK lexemes in current lexicographical descriptions, that is, the dictionaries *Suomen kielen perussanakirja* (PS) and *Nykysuomen sanakirja* (NS). As we can see in Table 5.21, 11 (6.1%) instances of features with significant odds in favor of a lexeme occurring were not exemplified at all in either dictionary. As a similar-sized discrepancy to the other direction, 12 (6.7%) features with significant odds against the occurrence of a particular lexeme were nevertheless portrayed among the examples. The discrepant feature-lexeme pairings in question are presented in Table 5.22 further below, in which it becomes evident that *ajatella* in particular is presented in the dictionaries in contexts which on the basis of the multivariate corpus analysis would be considerably more typical in conjunction with some other(s) of the selected THINK lexemes in terms of their odds. Such features attributed to *ajatella* in the dictionaries are ACTIVITY as a PATIENT, which ranks highest in terms of feature odds for *harkita*, as well as INDIRECT QUESTIONS as PATIENT which are likewise among the topmost ranked features for both *miettiä* and *pohtia*. Furthermore, whereas all features with significant odds in favor of the occurrence of *miettiä* are also represented in the dictionaries, there are in the case of *pohtia* no features among its example sentences which would have received significant odds against their co-occurrence with this particular lexeme. Thus, also the

multivariate results suggest that some of the examples presented in the dictionaries may not be the most characteristic ones for the studied four THINK lexemes, as much acceptable and possible that they might otherwise be.

Table 5.21. Comparison of lexeme-specific occurrences of features in the example sentences in the two dictionaries (PS and NS) against the directions of the odds either in favor of, against or neutral with respect to each lexeme derived with the multivariate polytomous logistic regression analysis of the research corpus; only features included in the proper full model (VI) are considered.

Dictionaries/Multivariate results	Significant odds in favor (>1)	Nonsignificant (neutral) odds (≈1)	Significant odds against (<1)	Σ
PS	22	31	5	58
NS	31	48	12	91
PS+NS	21	30	5	56
Ø	11	52	24	87
Σ	43	101	36	180

Table 5.22. Features with significant odds in favor of a lexeme but not exemplified in either dictionary (PS or NS), as well as features with significant odds against a lexeme but nonetheless exhibited among the dictionary example sentences.

Lexeme/Discrepancy	Features with significant odds in favor of a lexeme missing from both dictionaries (11)	Features with significant odds against a lexeme exemplified in either dictionary (12)
ajatella	PATIENT+PARTICIPLE (5.3:1) MANNER+GENERIC (23:1) MANNER+NEGATIVE (4:1)	AGENT+GROUP (1:5) PATIENT+NOTION (1:4.1) PATIENT+ACTIVITY (1:7.1) PATIENT+INDIRECT_Q... (1:14) MANNER+JOINT (1:2.7) TMP+INDEFINITE (1:1.7) DURATION (1:8.4) VERB_CHAIN+NECESSITY (1:2.9) CO-ORDINATED_VERB (1:2.1)
miettiä	-	PATIENT+että 'that' (1:1.9)
pohtia	PLURAL (1.6:1) AGENT+GROUP (4.2:1) PATIENT+INDIRECT_Q... (2.8:1) PATIENT+DIRECT_Q... (8.1:1) TMP+DEFINITE (2.3:1) VERB_CHAIN+TEMPORAL (2.4:1)	-
harkita	CONDITIONAL (2.3:1) META (1.6:1)	GOAL (1:4.7) QUANTITY (1:3)

Finally, we can make a small excursion to compare the results of the chosen proper full model (VI) with the additional descriptive intricacy allowed in the extended full model (VIII), in which the observed semantic subtypes are also included for the less frequent syntactic arguments. The complete results with all features and lexeme-specific odds are presented in Table R.19 in Appendix R, and while the greater size of the feature set has led to some individual changes throughout the entire set of feature-wise lexeme-specific odds, I will concentrate here on the features left out of the proper full model, that is, DURATION, FREQUENCY, LOCATION, QUANTITY, the more specific subtypes of modality for the verb chain, as well as the CO-ORDINATED VERBS.

In the case of DURATION, while the overall significant odds against *ajatella* persist in all subtypes, the significant odds in favor of *miettiinä* remain only for the LONG and SHORT but not for the OPEN and OTHER (referring to indication of a fixed temporal beginning or end point, or both) subtypes. For FREQUENCY, while the lumped OTHER (NON-OFTEN number of times) subtype is neutral for all the four lexemes, the overall significant odds against *ajatella* continue only with respect to the OFTEN but not the AGAIN subtypes, whereas the significant odds in favor of a lexeme are split between *miettiinä* for the OFTEN subtype and *harkita* for the AGAIN subtype. In the case of LOCATION, the overall significant odds against *ajatella* remain for all subtypes, but disappear in the case of *harkita*. The lexeme-specific significant odds in favor of *pohtia* persist for the (physical) LOCATION and EVENT subtypes, but are shifted to *miettiinä* in the case of the GROUP subtype. For QUANTITY, whereas the LITTLE subtype exhibits no significant lexeme-specific odds, the overall significant odds in favor of *miettiinä* apply for the MUCH subtype, but the general significant odds against *harkita* evaporate.

As for the semantic characterizations of the entire verb-chain, the overall neutrality of POSSIBILITY turns into significant odds against *pohtia* in the case of (positive) PROPOSSIBILITY and against *harkita* in the case of IMPOSSIBILITY. For NECESSITY, while NONNECESSITY is neutral for all four THINK lexemes, the overall significant odds in favor of *miettiinä* apply also to PRONECESSITY and FUTILITY, whereas in addition to the persistence of the overall significant odds against *ajatella* for both of these subtypes, *harkita* also becomes dispreferred in the case of PROPOSSIBILITY. Finally, for CO-ORDINATED VERBS as arguments, their significant odds in favor of both *ajatella* and *miettiinä* persist in the MENTAL but not the ACTION subtypes.

In conclusion, we can note that even though the semantic subtypes of these rarer arguments sometimes follow the preference patterns of the syntactic argument, at other times this is not the case, with preferences and dispreferences split among several lexemes or turning altogether neutral. Thus, when grouping semantic subtypes together we always lose some of the information contained in the research corpus. Nevertheless, this action must occasionally be taken in order to keep the number of feature variables within the recommended ratio with respect to the outcome frequencies in the research data set, as is the case in this study.

5.4 Assessing the robustness of the effects

5.4.1 Simple bootstrap and writer-cluster resampling

We can next move on to evaluate the robustness of the above observed effects, represented by the odds assigned to the explanatory variables, by applying two resampling procedures following the example presented in Section 3.4.4, namely, a 1000-fold simple bootstrap resampling already referred to in terms of the overall performance of the full model above, and a 10000-fold bootstrap with resampling from speakers or writers treated as clusters. With these magnitudes of iterations we can use the percentile method for calculating the confidence intervals of the statistics of our interest; the even greater number of repetitions for the speaker/writer-cluster scheme is motivated by the substantially smaller size of training data set which it by design uses, being in this study determined as 571 on the basis of the number of code-identified writers, whereas the simple bootstrap incorporates on the average 2152 ($\approx 3404 \cdot [1 - e^{-1}]$) unique (but each time random) instances of the altogether 3404 observations in the data (and some of such instances more than once).

While the simple bootstrap overfits to the training data only slightly, with a corresponding mean $R_L^2(TRAIN)=0.327$ and a 95% Confidence Interval of (0.307, 0.342), it is able to predict correctly outcomes in the entire data fairly well with a *Recall* of 63.8% (63.07-64.51%), $\lambda_{prediction}=0.355$ (0.343, 0.368), and $\tau_{classification}=0.479$ (0.468, 0.489), and its fit with this testing data is also acceptable with a mean $R_L^2(TEST)=0.287$ (0.264, 0.300). In contrast, whereas the prediction accuracy of the model trained with writer-cluster bootstrap resamples does not fall drastically in comparison to the simple bootstrap, having for the entire testing data a mean *Recall*=60.6% with a 95% Confidence Interval of (58.93-62.07%), $\lambda_{prediction}=0.299$ (0.269, 0.325), and $\tau_{classification}=0.433$ (0.408, 0.454), it overfits considerably more with the smaller training data, with $R_L^2(TRAIN)=0.400$ (0.357, 0.446). Despite the relatively good success in outcome prediction, this overfit results on the average in an utterly dismal fit with the testing data, as mean $R_L^2(TEST)=-0.118$ (-0.375, 0.067).

The negative R_L^2 values, which mean that the trained model performs on the testing data worse than the null model using relative frequencies of the lexemes as its default estimated probabilities, results from the extremization of outcome probability estimates due to the overfit at the training stage, which is reflected and mediated by more extreme parameters (i.e., logarithms of the odds) assigned to the lexeme-wise feature variables in the model. As was noted in Section 3.4.4, only a proportionately small number of original outcomes assigned with a relatively small probability by a fitted model can readily increase the model deviance D_{model} over the null Deviance D_{null} , resulting in negative R_L^2 values. In light of these results, the consideration of some of other, more sophisticated measures of model fit presented in statistical literature (see, e.g., Mittlböck and Schemper 1996, 2002; also Hosmer and Lemeshow 2000: 144-156, 164-167) would seem recommendable in future studies.

Table 5.23. The odds assigned to the two semantic subtypes of AGENT for each lexeme in the proper full model (VI) with a single fit from the entire data ($n=3404$) using the one-vs-rest heuristic; odds against any lexeme, i.e., $\beta(L|F) < 0 \sim e^{\beta(L|F)} < 1$, supplemented by the corresponding ratio, i.e., $1:1/e^{\beta(L|F)} \sim 1:e^{-\beta(L|F)}$, e.g., 0.5~1:2; significant lexeme-wise odds in **boldface**; nonsignificant odds in (parentheses); features with at least one lexeme with significant odds in **boldface**.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected\$odds.mean

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX AGE.SEM GROUP	0.2~1:5	0.52~1:1.9	4.2	(1.1)
SX AGE.SEM INDIVIDUAL	(0.85~1:1.2)	(0.98~1:1)	(1.6)	(0.69~1:1.5)

Table 5.24. The 95% Confidence Intervals for the odds assigned to the two semantic subtypes of AGENT for each lexeme in the proper full model (VI) using 1000-fold simple bootstrap resampling from the entire data ($n=3404$); results differing from the original single-round fit with the entire data marked with thicker border-lines, such odds having turned from nonsignificant to significant *italicized*.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected.1000\$odds.range

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX AGE.SEM GROUP	0.12<..<0.29	0.25<..<0.94	2.6<..<6.7	(0.54<.. <2.1)
SX AGE.SEM INDIVIDUAL	(0.44<.. <1.2)	(0.61<.. <1.7)	1.03<..<2.9	(0.35<.. <1.3)

Table 5.25. The 95% Confidence Intervals for the odds assigned to the two semantic subtypes of AGENT for each lexeme in the proper full model (VI) using 10000-fold bootstrap resampling from writers ($n=571$) as clusters; results differing from the original single-round fit with the entire data marked with thicker border-lines, such odds having turned from significant to nonsignificant ~~struck through~~.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected.10000_speaker_cluster\$odds.range

Feature/Lexeme	ajatella	miettiä	pohtia	harkita
SX AGE.SEM GROUP	0.024<..<0.57	(0.096<..<2.3)	(0.82<..<9.1)	(0.2<.. <5.3)
SX AGE.SEM INDIVIDUAL	(0.35<.. <3.1)	(0.29<.. <4.2)	(0.43<.. <4.2)	(0.046<.. <1.8)

Nevertheless, following Bresnan et al. (2007) my focus is rather to assess the robustness of the explanatory features in the proper full model indicated by the range of their estimated odds over the resamples in the two schemes. The full results in this respect are presented in Tables R.20-R.21 in Appendix R, of which an exemplary sample is shown above in Tables 5.23-5.25 concerning the two semantic subtypes of AGENT among the four selected THINK lexemes. As can be seen by comparing Table 5.23 with the single-fit odds against 95% Confidence Intervals derived with the 1000-fold simple bootstrap in Table 5.24, the range 2.6..6.7 corresponding to *pohtia* in conjunction with a GROUP as AGENT both encompasses the single-fit significant odds 4.2, and stays clearly above 1.0 thus validating the significance and its direction for this feature in favor of *pohtia*.

Likewise, the ranges 0.12..0.29 for *ajatella* and 0.25..0.94 for *miettiä* in conjunction with the same feature lie below 1.0 and encompass the corresponding significant single-fit odds of 0.2 and 0.52, respectively, against the occurrence of these lexemes, though *miettiä* does come relatively close, but not quite, to non-significance, that is, bridging both sides of 1.0. The latter is the case with the range 0.52..2.1 for *harkita*, still in conjunction with the same GROUP as AGENT, which reaffirms the non-

significance of the corresponding single-fit odds of 1.1, and the same situation applies also for the odds and odds-ranges of *ajatella*, *miettiinä*, and *harkita* in conjunction with an INDIVIDUAL as AGENT.

In contrast, the range 1.03..2.9 for *pohtia* in conjunction with the latter feature suggests that the assessment of the corresponding single-fit odds of 1.6 as nonsignificant rather than in favor of the lexeme in question is a border-line case. Indeed, as can be seen from the overall statistics in Table 5.26 below, the 1000-fold simple bootstrap renders as significant altogether 96 combinations of features in conjunction with the studied THINK lexemes, which is 13 more than the result from the single fit, and this increase is lexeme-wise evident for all lexemes but *ajatella*. In all, this amounts to two more features with at least one significant lexeme-specific odds, concerning the already mentioned INDIVIDUAL as AGENT as well as CLAUSE-EQUIVALENT forms, the latter which now turn significantly, though slightly, in favor of *harkita* and against *miettiinä*.

However, when we turn to the 95% Confidence Intervals derived using the 10000-fold bootstrap with resampling from writers treated as clusters, we can see in Table 5.25 that the only significant lexeme-specific odds range which remains is assigned for *ajatella* in conjunction with a GROUP as AGENT, being 0.024..0.57 and consequently against their co-occurrence, while the ranges for all the other cases expand to extend to both sides of 1.0 and thus indicate non-significance. Overall, the number of significant combinations of features with the studied lexemes falls down drastically to 38, less than half the corresponding figures for both the single fit and the 1000-fold simple bootstrap, and this reduction applies to all four THINK lexemes. Consequently, also the number of features with at least one significant lexeme-specific odds drops to 20. This shows that at least a part of the effects observed in the single-fit odds are not strong enough to not be *possibly* attributable to individual writer/speaker preferences, that is, they are not sufficiently dispersed and frequent among the entire writer/speaker population to remain significant when individual speakers are randomly sampled instead of all the individual usage instances.

On the other hand, the said 20 features which continue to exhibit significant odds with respect to the studied THINK lexemes, despite this harsher sampling scheme, can with justification be concluded to be writer-independent, and thus the most robust features incorporated in the proper full model. The mean odds resulting from the 10000-fold writer-cluster resampling for these 20 most robust features are presented in Table 5.27 below. As can be seen, the odds even as mean values are considerably more extreme than those in the single-fit model in Table 5.17 above, and this helps to explain how the estimated probabilities also become extreme as a function of the odds, resulting then in poorer fit when tested with the entire research data set, as was discussed earlier above. Nevertheless, such a poor fit does not diminish the apparent robustness of the features in question, although the actual odds values can be assumed to be all too extreme to be correct as such in describing the use of the studied THINK lexemes in general, outside the current data set.

Table 5.26. Results concerning the significance and confidence intervals of the odds estimated for the proper full model (VI) using both a 1000-fold simple bootstrap resampling procedure and a 10000-fold bootstrap with resampling from writers/speakers as clusters; odds for features with respect to a lexeme in the model are considered significant if their 95% Confidence Interval is fully either above or below 1.0; otherwise, the particular odds are considered nonsignificant; overall number of feature variables in the full proper model being 46, resulting in altogether 184 lexeme-feature pairings.

Significant odds/ Models Lexemes	Single-fit model	1000-fold simple bootstrap resampling	10000-fold bootstrap with resampling from writers as clusters
Features with at least one significant lexeme-specific odds	39	41	20
Features with all lexemes having significant odds	2	5	1
Features with no significant lexeme- specific odds	7	5	26
<i>ajatella</i>	31	31	18
<i>miittiä</i>	21	26	5
<i>pohtia</i>	19	22	8
<i>harkita</i>	12	17	7
Overall significant lexeme-specific odds	83	96	38
Overall nonsignificant lexeme-specific odds	101	88	146

Table 5.27. Mean odds of the 20 most robust feature variables selected from the proper full polytomous logistic regression model (VI) fitted using the one-vs-rest heuristic on 10000-fold bootstrap resamples from writers as clusters, with each of the studied THINK lexemes pitted against the others at a time; odds against any lexeme ($e^{\beta[L|F]} < 1$) supplemented by the corresponding ratio ($1/e^{\beta[L|F]} \sim e^{-\beta[L|F]}$, e.g., 0.5~1:2); significant lexeme-wise odds in **boldface**; nonsignificant odds in (parentheses); all features with at least one lexeme with significant odds.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_se
mantics.selected.10000.speaker.cluster\$odds.mean

Feature/Lexeme	ajatella	harkita	miettiä	pohtia
SX_AGE.SEM_GROUP	0.14~1/7.4	(1.1)	(0.54~1/1.8)	(2.9)
SX_DUR	0.011~1/94	(0.47~1/2.1)	4.2	(0.86~1/1.2)
SX_GOA	5.5	0~1/2.9e³	(0.2~1/5.1)	(0.14~1/7.1)
SX_LOC	0.2~1/4.9	(0.31~1/3.2)	(0.92~1/1.1)	4.1
SX_LX_että_CS.SX_PAT	3.6	(0.022~1/45)	(0.45~1/2.2)	(0.11~1/9)
SX_MAN.SEM_AGREE...	8.0e⁵	0~1/2.4e⁷	(0~1/4.2e ⁵)	(0~1/2.1e ⁶)
SX_MAN.SEM_GENERIC	4.6e⁵	0~1/3.7e⁷	0~1/1.8e ⁴	0~1/2.4e⁷
SX_PAT.DIRECT_QUOTE	0~1/3.4e³	0~1/2.4e⁷	(2.8)	8.3
SX_PAT.INDIRECT_Q...	0.053~1/19	(0.7~1/1.4)	3.8	3.6
SX_PAT.INFINITIVE	715	(0.001~1/1.3e ³)	0~1/2.2e⁷	(0~1/4.0e ⁴)
SX_PAT.PARTICIPLE	735	(0.005~1/210)	0~1/1.8e⁷	(0~1/6.8e ⁵)
SX_PAT.SEM_ABSTR...	0.25~1/4	(1)	(1.6)	3.7
SX_PAT.SEM_ACTIVITY	0.14~1/7.4	(9.8)	(0.71~1/1.4)	(1.6)
SX_PAT.SEM_COMM...	0~1/6832	(0.74~1/1.4)	(0.86~1/1.2)	(0.01~1/100)
SX_PAT.SEM_INDIV..._GR...	15	(0~1/1.6e ⁴)	(0.11~1/9)	0~1/2.1e⁴
SX_SOU	29	0~1/8.3e⁶	(0.01~1/102)	(0.001~1/1.8e ³)
SX_TMP.SEM_DEFINITE	(0.34~1/2.9)	(0.12~1/8.3)	(0.74~1/1.3)	3.1
SX_VCH.SEM_ACCIDENT...	(105)	0~1/3.5e⁷	(0~1/1.2e ⁴)	(0.007~1/150)
SX_VCH.SEM_NECESSITY	0.35~1/2.9	(1.5)	(1.9)	(0.91~1/1.1)
SX_VCH.SEM_TEMPORAL	0.066~1/15	(0~1/1.3e ⁵)	(1.7)	4.8

5.4.2 Assessing the effect of incorporating extralinguistic features

Lastly, we can scrutinize what type of impact adding extralinguistic features concerning the medium and context of language usage on top of the actual linguistic features included in the proper full model (VI) has on the odds, when estimated with a single fit using the entire data. The two extra-linguistic features firstly indicate whether an instance of the studied THINK lexemes has been used in newspaper text or in Internet newsgroup discussion (Z_EXTRA_SRC_sfnet), and secondly whether an instance within the newspaper subcorpus is part of a citation, typically representing a spoken fragment (Z_QUOTE). The specific results concerning the estimated odds for this selection of features are presented in full in Table R.22 in Appendix R, so I will note here in Table 5.28 below only the essential differences with respect to the plain model consisting only of proper linguistic feature variables.

In the first place, the two extralinguistic variables are both overall significant, in that the indication of the medium has significant odds with respect to all lexemes, as is also the case with the indication of usage within citation/quotation with all lexemes but *harkita*. Thus, having the Internet newsgroup discussion as the medium increases the odds in favor of both *ajatella* (1.6:1) and *miettiä* (2:1) occurring, while it decreases the odds of occurrence for both *pohtia* (1:2.2) and *harkita* (1:2.1). For its

part, usage within a citation increases the odds in favor of both *ajatella* (1.6:1) and *mieltä* (1.5:1), whereas it lowers those for *pohtia* (1:2), with *harkita* remaining nonsignificant. These results fit nicely with the hypotheses suggested in Section 4.1.2 that both *ajatella* and *mieltä* would be more partial than the others among the THINK lexemes to personal expression and direct discourse with identified recipients.

Overall, fitting the model including the extralinguistic characteristics results in 40 features with at least one significant lexeme-specific odds, of which 3 features have significant odds either way for all the lexemes, while 8 features receive no significant lexeme-specific odds. Lexeme-wise, there are altogether 33 features with significant odds for *ajatella*, 24 for *mieltä*, 16 for *pohtia*, and 12 for *harkita*. Though maximally the odds estimated for a linguistic feature can grow by a factor of 1.44:1 or diminish by a similar inverse factor, on the average the changes are quite small, being approximately 1.08:1 or its inverse. The greatest increases of the lexeme-specific odds concern DIRECT QUOTES as PATIENT with both *ajatella* and *mieltä*, and CONDITIONAL mood with *harkita*, whereas the greatest corresponding decreases involve a GROUP as AGENT with *pohtia*, a DIRECT QUOTE as PATIENT with *mieltä*, ACTIVITY as PATIENT with *harkita*, and SECOND person with *mieltä*. In no circumstances do significant odds in favor of a lexeme in conjunction with a feature turn into significant odds against the same lexeme, or *vice versa*, when the two extralinguistic features are incorporated in the model. However, significant odds may turn nonsignificant, as is the case with AGREEMENT and NEGATIVE evaluation as subtypes MANNER as well as ACTIVITY as PATIENT and PASSIVE voice manifested in the verb chain with *pohtia*. The last-mentioned change also entails that the PASSIVE voice becomes nonsignificant overall with respect to all four THINK lexemes. In contrast, the opposite change from nonsignificant to significant odds happens in only two cases, namely, with NEGATION in conjunction with both *mieltä* and *pohtia*. Nevertheless, the overall impact of the extralinguistic features does not seem to be that substantial, at least when considered as such, without calculating their interactions with the actual linguistic variables (as is tentatively explored in Appendix M).

Table 5.28. Selected odds of the proper full polytomous logistic regression model supplemented with the two extra-linguistic variables (Model VIII), fitted using the one-vs-rest heuristic from the entire data ($n=3404$), with each of the studied THINK lexemes pitted against the others at a time; nonsignificant odds in (parentheses); odds against any lexeme ($x<1$) supplemented by the corresponding ratio ($1/x$, e.g., $0.5\sim 1:2$); significant lexeme-wise odds in **boldface**; nonsignificant odds in (parentheses); features with at least one lexeme with significant odds in **boldface**, results differing from the original single-round fit with the entire data with thicker border-lines, such odds having turned from significant to nonsignificant ~~struck through~~, those from nonsignificant to significant *italicized*; significant odds which have changed by more than the mean difference marked with ‘*’.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected.extra\$odds.mean

Feature/Lexeme	ajatella	miettä	pohtia	harkita
SX MAN.SEM AGREEMENT	16	0.068~1:15	(0.24~1:4.2)	(0~1:6.7e ⁶)
SX MAN.SEM NEGATIVE	3.9	(0.53~1:1.9)	(0.24~1:4.2)	(0.65~1:1.5)
*SX PAT.SEM ACTIVITY	*0.15~1:6.5	(0.9~1:1.1)	(1.3)	*7.7
*Z ANL NEG	2	0.68~1:1.5	*0.53~1:1.9	(1.2)
Z ANL PASS	(0.66~1:1.5)	(1)	(1.7)	(0.97~1:1)
Z EXTRA_SRC sfnet	1.6	2	0.45~1:2.2	0.47~1:2.1
Z QUOTE	1.6	1.5	0.49~1:2	(0.91~1:1.1)

We can now also reassess the average weights of the various feature categories when the extralinguistic features are incorporated. As can be seen in Table 5.29, the extralinguistic features have the lowest mean aggregate odds (1.86), bettering morphological features manifested in the node-verb or the enveloping verb chain only in the special case of mean odds *against* the occurrence of a lexeme (1:2.13 vs 1:1.91) and in the aggregate case when also nonsignificant odds are considered (1.74 vs 1.46). With respect to the other feature categories, syntactic arguments coupled with their semantic and structural subtypes remain as the most distinctive group with a mean aggregate odds of 4.13:1, whether in favor (3.71) or against (1:4.70) the occurrence of a lexeme. At some distance, the semantic characterizations of the verb chain reach second place, followed closely by syntactic arguments by themselves, with mean aggregate odds of 3.17:1 and 2.92:1, respectively. However, when considering only odds in favor of lexemes, the syntactic arguments fare slightly better than verb chain general semantic characterizations, the odds being 2.69:1 vs. 2.66:1. Comparing these results overall with those presented in Table 5.16 in Section 5.3, we can conclude that the particular extra-linguistic features have practically no bearing on the absolute weights and relative ranking of the other feature categories, which is probably also reflected in that extra-linguistic features are relegated to the lowest rank relative to the rest with respect to the magnitude of their lexeme-specific odds (cf. also the lowest performance figures in Table 5.14 in Section 5.2.2 for Model X consisting of only the two extra-linguistic variables).

Table 5.29. Average weights of the different categories of explanatory feature variables including the two extralinguistic ones, calculated firstly for significant odds in the overall polytomous regression model (VIII) attained with the one-vs-rest heuristic on the entire data, and secondly for all estimated odds including nonsignificant ones (in parentheses).

Feature variable category	Mean odds in favor	Mean odds against	Mean aggregate odds
Verb chain morphology	2.02 (1.48)	0.52~1:1.91 (0.70~1:1.44)	1.96 (1.46)
Verb chain semantics	2.66 (1.62)	0.24~1:4.24 (0.12~1:8.59)	3.17 (3.48)
Syntactic argument types (alone)	2.69 (1.84)	0.32~1:3.14 (0.47~1:2.11)	2.92 (1.99)
Syntax arguments + semantic/structural subtypes	3.71 (2.57)	0.21~1:4.70 (0.06~1:18)	4.13 (7.89)
Extralinguistic features	1.68 (1.68)	0.47~1:2.13 (0.56~1:1.80)	1.86 (1.74)

5.5 Probability estimates of the studied lexemes in the original data

5.5.1 Overall assessment of lexeme-wise probability estimates

In addition to assigning for explanatory variables parameter values which can be interpreted as odds, as discussed at length above, the other attractive characteristic of a (polytomous) logistic regression model is its ability to provide probability estimates for an outcome, given any possible mix of explanatory variables, representing a set of features present in some context. Like the estimated odds, the accuracy of such probability estimates is naturally dependent on how well the explanatory variables incorporated in the model are able to describe and fit the data they are trained with, as well as to predict instances in new, unseen data, that is, how generally applicable the selected model is. Nevertheless, the probability estimates allow us to effectively rank with a single value the joint effect of a large number of features and their complex interrelationships, which is typically the case with real, natural usage of language.

Reminiscing the sentences in the original data containing instances of the studied THINK lexemes and the practical reality of conducting their linguistic analysis, there are very few clean and clear cases, where one could easily isolate only one or two significant feature variables and consider the rest as neutral or altogether ignorable. Natural language usage is difficult if not impossible to reduce to simple “laboratory sentences”, and such artificially constructed combinations of thoroughly controlled linguistic items and nothing else are lacking in naturalness in the eyes or ears of a (native) language user. With a logistic regression model we can systematically rate entire sentences, or even longer text fragments, together with all the relevant linguistic information they contain with respect to the studied linguistic phenomenon. In estimating the probability ratings to be scrutinized in depth below, I have decided to include also the two extralinguistic variables in addition to the proper linguistic ones, corresponding to Model VIII as described in Section 5.1. Though the overall impact of extralinguistic variables appears to be relatively low in comparison to the linguistic ones, they can nevertheless be considered relevant as proxies for the style of linguistic usage, as either impersonal, detached narration/reporting or personal, immediate discourse, as well as the expected level of adherence to linguistic norms, as either formal or informal.

I will first look into the sums of the probabilities estimated for the individual lexemes, since as the associated models are fit separate of each other they do not necessarily add up to the theoretically correct $\sum_{Lexeme} P(Lexeme|Context)=1.0$. As we can see in Figure 5.4 below, there is clearly dispersion in the sum probabilities, ranging at the extremes from a minimum of $\sum P=0.546$ to a maximum of $\sum P=1.711$. However, the mean of the sum probabilities is clearly 1.0, around which the bulk of the values are tightly concentrated, as the 95% Confidence Interval is already $CI(\sum P)=(0.771, 1.195)$, excluding the outlier. This span roughly coincides with what one could conclude by visual inspection of Figure 5.4. Nevertheless, tightening the Confidence Interval further narrows the interval of sum probabilities down only gradually, as the 90% range $CI(\sum P)=(0.826, 1.139)$, the 80% range $CI(\sum P)=(0.878, 1.102)$, and the 50% range still $CI(\sum P)=(0.944, 1.057)$.

Consequently, one can conclude that the separately fit individual lexeme-specific binary models constituting the overall polytomous model do not produce a perfect fit,

but nevertheless they can be considered to roughly approximate the ideal case. Thus, in the following scrutinies the lexeme-specific probabilities for each instance in the data set are adjusted so that their sum per each such instance will equal $\sum P=1.0$. This is done by simply dividing instance-wise each original lexeme-specific probability estimate by the sum of these estimates for that particular instance, i.e., $P_{adjusted}(Lexeme|Context)=P_{original}(Lexeme|Context)/\sum P_{original}(Lexeme|Context)$. These adjusted probability values were already used in the calculation of model fit with the training and the testing data using the R_L^2 statistic, which is through model deviances based on individual instance-specific probability estimates.

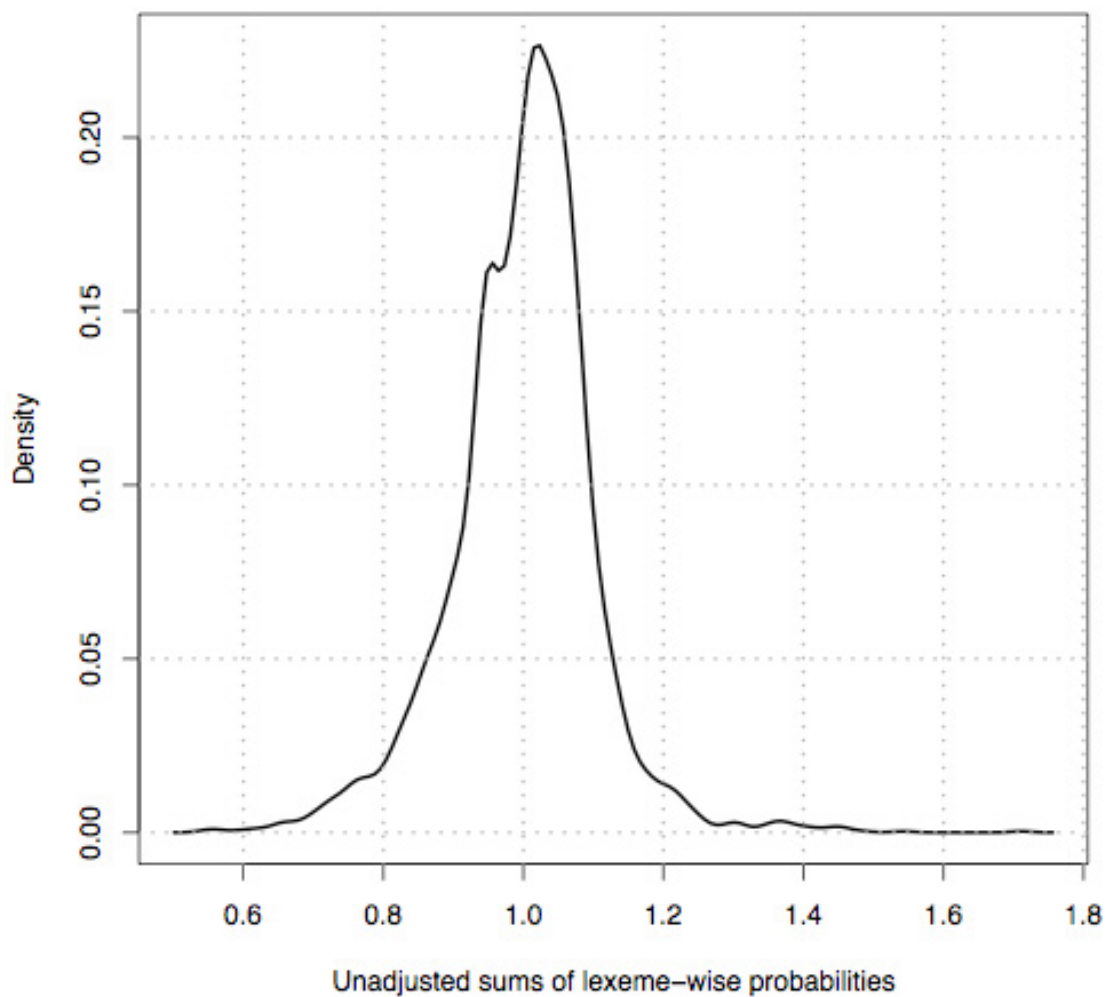


Figure 5.4. The distribution of the unadjusted sums of lexeme-wise probabilities on the basis of proper full model plus extralinguistic features (VIII) for each instance in the entire data set ($n=3404$).

Next, the underlying premises of logistic regression analysis, that is, assuming relative proportions of occurrence rather than categorical selections, suggest that we look not only at the maximum probabilities assigned for each instance but the entire spectrum of probabilities estimated for each outcome (i.e., $Lexeme \sim L$) in a particular context ($\sim C$). Indeed, as we can see in Figure 5.5, the maximum probability assigned for any

lexeme in any context rarely approaches the theoretical maximum $P(L|C)=1.0$, and the predictions are practically categorical in only 258 (7.6%) instances for which $P_{max}(L|C)>0.90$. To the contrary, the mean maximum probability per all instances and contexts is only $\bar{x}(P_{max}[L|C])=0.636$, while the overall span of maximal values is as broad as (0.28, 1.00), and even the 95% CI=(0.369, 0.966). The lower-ranked instance-wise probability estimates have similar overall characteristics of intermediate-level means and broad ranges. The second-highest probability estimates per instances have a mean $\bar{x}(P_{max-1}[L|C])=0.244$, with an overall range of (0.000, 0.490) and a 95% CI=(0.026, 0.415), and the third-highest probability estimates have a mean $\bar{x}(P_{max-2}[L|C])=0.096$, with an overall range of (0.000, 0.307) and a 95% CI=(0.000, 0.241). Even the minimum probability estimates clearly keep some distance from zero as their mean $\bar{x}(P_{min}[L|C])=0.043$, even though their overall range is (0.000, 0.212) as well as 95% CI=(0.000, 0.144). Nevertheless, as many as 764 (22.4%) of the minimum estimated probabilities per instance are practically nil with $P_{min}(L|C)<0.01$. However, turning this the other way around, for 2640 (77.6%) instances the minimum estimated probability is $P_{min}(L|C)\geq 0.01$, that is, representing an expected possibility of occurrence at least once every hundred times or even more often for *all four* THINK lexemes in a similar context.

Looking at the instance-wise estimated probabilities as a whole, in 64 (1.9%) instances all four estimates are $P\geq 0.15$, indicating relatively equal values for all lexemes, and in 331 (9.7%) instances all four are $P\geq 0.10$. Discarding always the minimum value, in 303 (8.9%) cases the remaining three higher-ranked probability estimates are all $P\geq 0.2$, and in as many as 1436 (42.2%) cases $P\geq 0.10$. Narrowing our focus only to the two topmost-ranked lexemes per instance, in 961 (26.2%) cases both probability estimates are $P\geq 0.3$, and for as many as 150 (4.4%) cases both $P\geq 0.4$. The contextual settings associated with these last-mentioned instances would be prime candidates for fully or partially synonymous usage within the selected set of THINK lexemes, as their joint probabilities would indicate high mutual interexchangeability. In sum, these distributions of instance-wise probability estimates for all four THINK lexemes suggest that, to the extent these probabilities even approximately represent the proportions of actual occurrences in the given contexts, very few combinations of contextual features are associated with categorical, exception-less outcomes. On the contrary, quite many of the contexts can realistically have two or even more outcomes, though preferential differences among the lexemes remain to varying extents (cf. Hanks 1996: 79). In that the contextual features used in this study are good and satisfactory representatives of a theory of language, that is, the fundamental components of which language is considered to consist and with which language can be comprehensively analyzed, as well as the rules or regularities concerning how these component parts interact and are allowed to combine in sequences, these results certainly support Bresnan's (2007) probabilistic view of the relationship between language usage and the underlying linguistic system. As we shall see in Section 5.5.2, the instance-wise context-based probability estimates are not merely an artefact resulting from applying a probabilistic method to the data, but correspond to actual proportions evident in the data (which logistic regression specifically tries to model).⁹³

⁹³ I am grateful to my preliminary examiner Stefan Th. Gries for drawing my attention to the problematics of this probabilistic view.

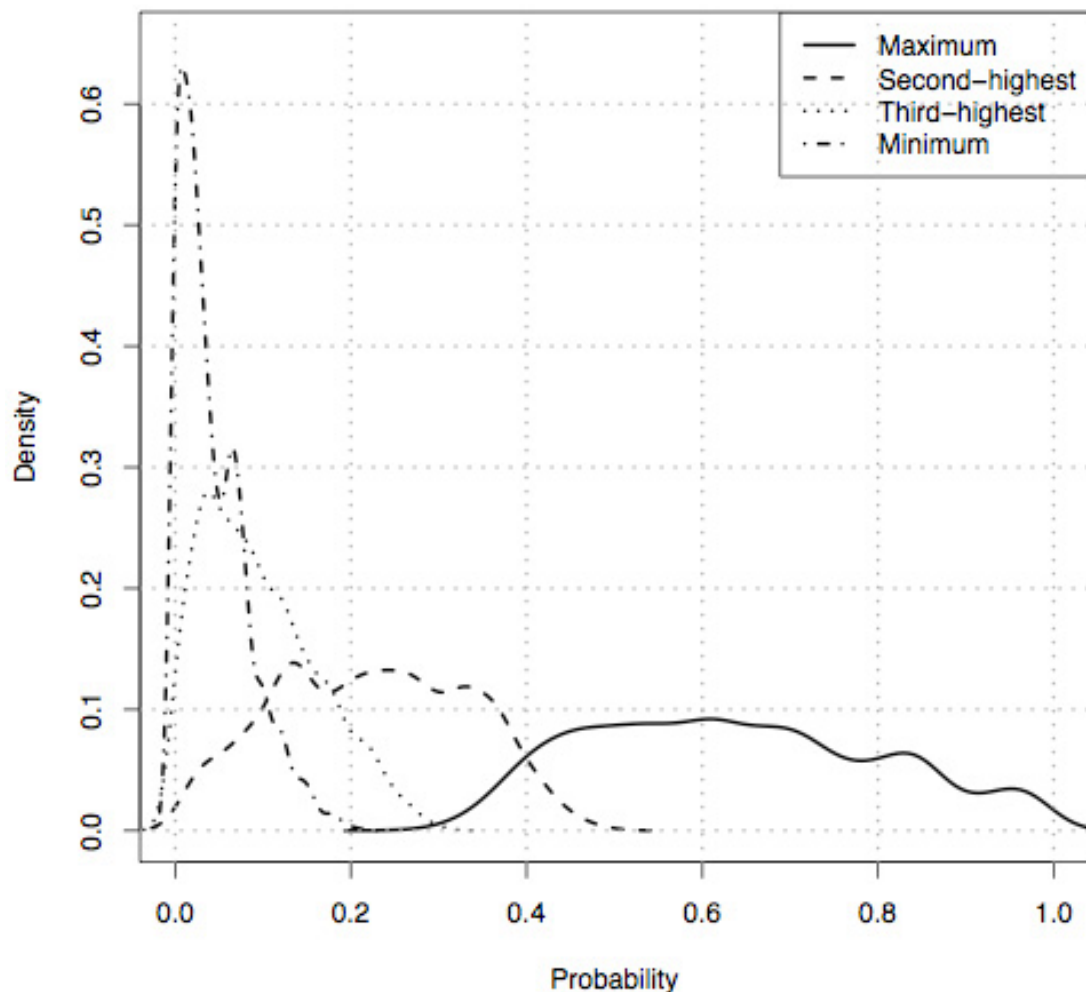


Figure 5.5. Densities of the distributions of the estimated probabilities by ranking order for all instances in the data ($n=3404$).

Next, we can turn to the lexeme-specific probability estimates, presented for all four THINK lexemes over the entire data in Figure 5.6. It should not be surprising that all lexeme-specific distributions are skewed towards the left with lower probability, albeit to different degrees, with the heights of their maximal peaks at that end corresponding inversely with their overall frequency in the data. Consequently, while for *ajatella* the mean probability is still close to the center with $P=0.437$, with an overall maximal range of (0.000, 1.000), and 95% $CI=(0.011, 0.966)$, for the rarer lexemes these general values are lower and the spans narrower, as for *miettiä* the mean is $P=0.241$, the overall range (0.000, 0.889), and the 95% $CI=(0.00, 0.73)$, for *pohtia* the mean is $P=0.210$, the range (0.000, 0.852), and the 95% $CI=(0.00, 0.69)$, and finally for *harkita* the mean is $P=0.113$, the overall range (0.000, 0.725), and the 95% $CI=(0.000, 0.558)$.

Focusing on the peaks and contours in the lexeme-specific probability distributions, we can see that the lexemes have different numbers of local maxima with varying positions on the probability range. The most frequent lexeme, *ajatella*, has also the

most even spread, with one level top broadly around $P \approx 0.8$ and another round one just below $P \approx 0.1$, thus towards both ends of the probability range, suggesting the underlying binary model has a propensity to either give relatively strong probabilities either in favor of or against the occurrence of this lexeme. In turn, the probabilities for *miettiä* climb slowly to a plateau towards around $P \approx 0.6$, rising then to one round hump just below $P \approx 0.4$ and a higher one at $P \approx 0.1$, staying mostly above the two other less frequent lexemes until the very lowest probability range. For its part, *pohtia* has two gentle tops at around $P \approx 0.6$ and $P \approx 0.4$ before the highest one somewhat below $P \approx 0.1$. Finally, *harkita* has two little upward bumps on both sides of $P \approx 0.5$ and a third one just below $P \approx 0.3$, before rising sharply to the highest peak of all at barely above $P \approx 0.0$. These peaks among the probability distributions for the three less frequent THINK lexemes suggest that for each there are a few specific contextual feature combinations which are particularly frequent in comparison to the other evident contexts. It is also possible that such clustering of probabilities around a small set of values may to some extent arise from the properties of the mathematical process by which the polytomous logistic regression model is fit with the data.

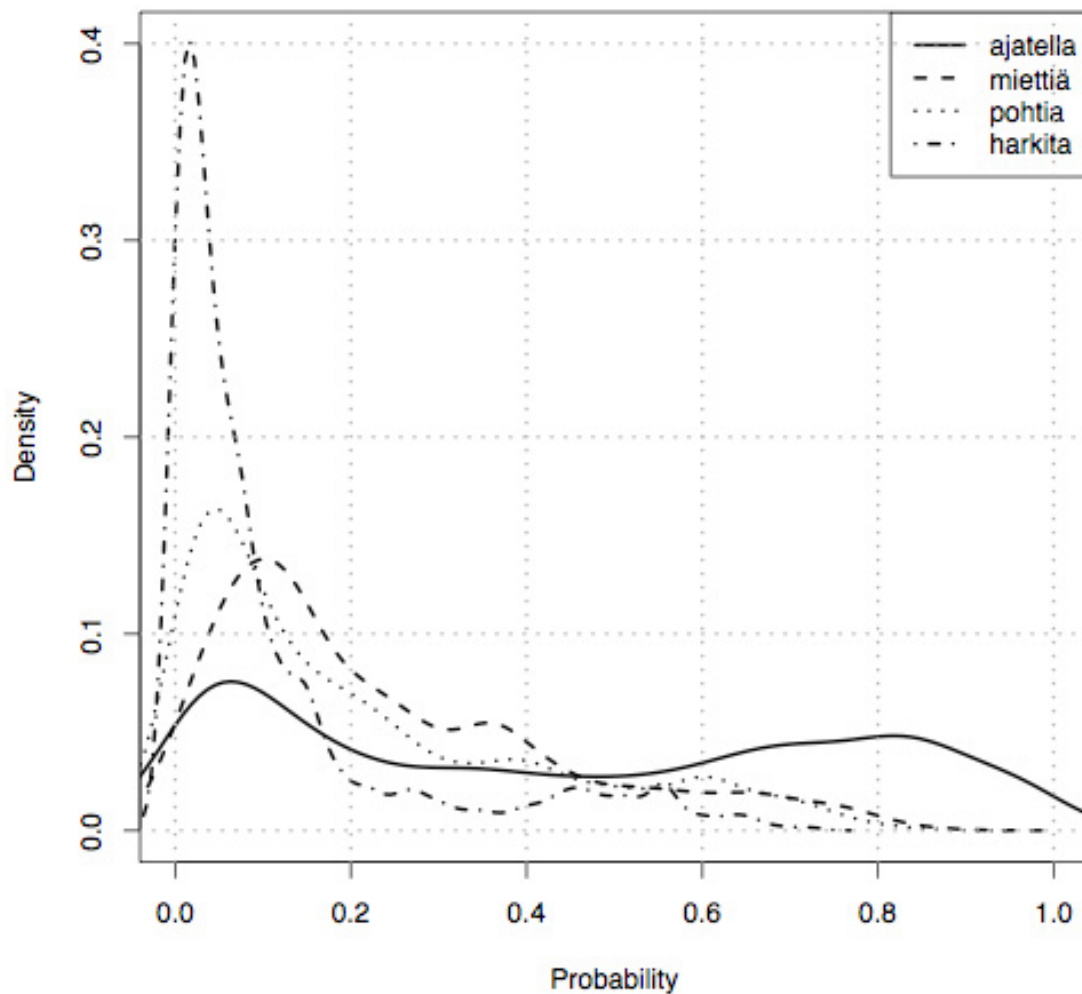


Figure 5.6. Densities of the distributions of the estimated probabilities by each lexeme for all instances in the data ($n=3404$).

These probability peaks are rendered more prominent when we look at the distributions of probability estimates per lexeme for those cases in which these have been the highest among all four lexemes for some instance in the data set, that is, when the prediction rule $\arg_{Lexeme}\{\max[P(Lexeme|Context)]\}$ would suggest the selection of the lexeme in question in that particular instance with its composition of feature variables, regardless of whether this selection would actually match the original lexeme in the corpus data or not. These distributions in Figure 5.7 show that *ajatella* would be the selection of choice with on the average highest estimated probabilities, typically $P > 0.6$, whereas the bulk of probability estimates with which either *miettä* or *pohtia* would be selected lies in the intermediate range roughly between $P \approx (0.4, 0.8)$, while *harkita* would be selected with the relatively weakest probability estimates between $P \approx (0.4, 0.6)$. In particular, *ajatella* has its peak at being selected at just above $P \approx 0.8$, *pohtia* at approximately $P \approx 0.6$, and *harkita* at two peaks just on each side of $P \approx 0.5$, while *miettä* has a relatively flat and broad maximum plateau between $P \approx (0.4, 0.5)$.

The distributions also indicate that among the three rarer lexemes *miettiinä* and *pohtia* are seldom assigned the very highest possible probabilities, as the maxima discernible already in the overall value ranges above are $P=0.889$ for *miettiinä* and $P=0.852$ for *pohtia*, and even lower for *harkita* at $P=0.725$. Thus, the combination of the underlying binary logistic models constituting the polytomous model clearly exhibits its highest confidence in the prediction of *ajatella*, the most frequent among the lot, whereas the predictions for the three rarer THINK lexemes mostly leave some room, in the form of substantial “leftover” probability, for one or more of the other lexemes to possibly occur. This can be also understood to entail that such contexts might exhibit genuine, permissible variation among the studied lexemes, in which case at least the current linguistic variables, and possibly any extension of such set, would not produce a categorical choice.

Furthermore, the maximal values of the probability estimates for the three rarer lexemes are clearly less than those preliminarily presented in Arppe (2007), but this is due to their normalization here to $\sum P=1.0$, rather than any substantially improved effects from the somewhat larger feature set employed in that study, corresponding in composition to approximately the extended full model (VIII) in this dissertation. In fact, if we look at the unadjusted original possibilities for the rarer THINK lexemes, their maxima are certainly higher, being $P=0.957$ for *miettiinä*, $P=0.911$ for *pohtia*, and $P=0.914$ for *harkita*. Thus, adjusting the probabilities to adhere to $\sum P=1.0$ does somewhat penalize the estimates for the rarer THINK lexemes, as these tend to be, on the whole, lower than those for *ajatella*.

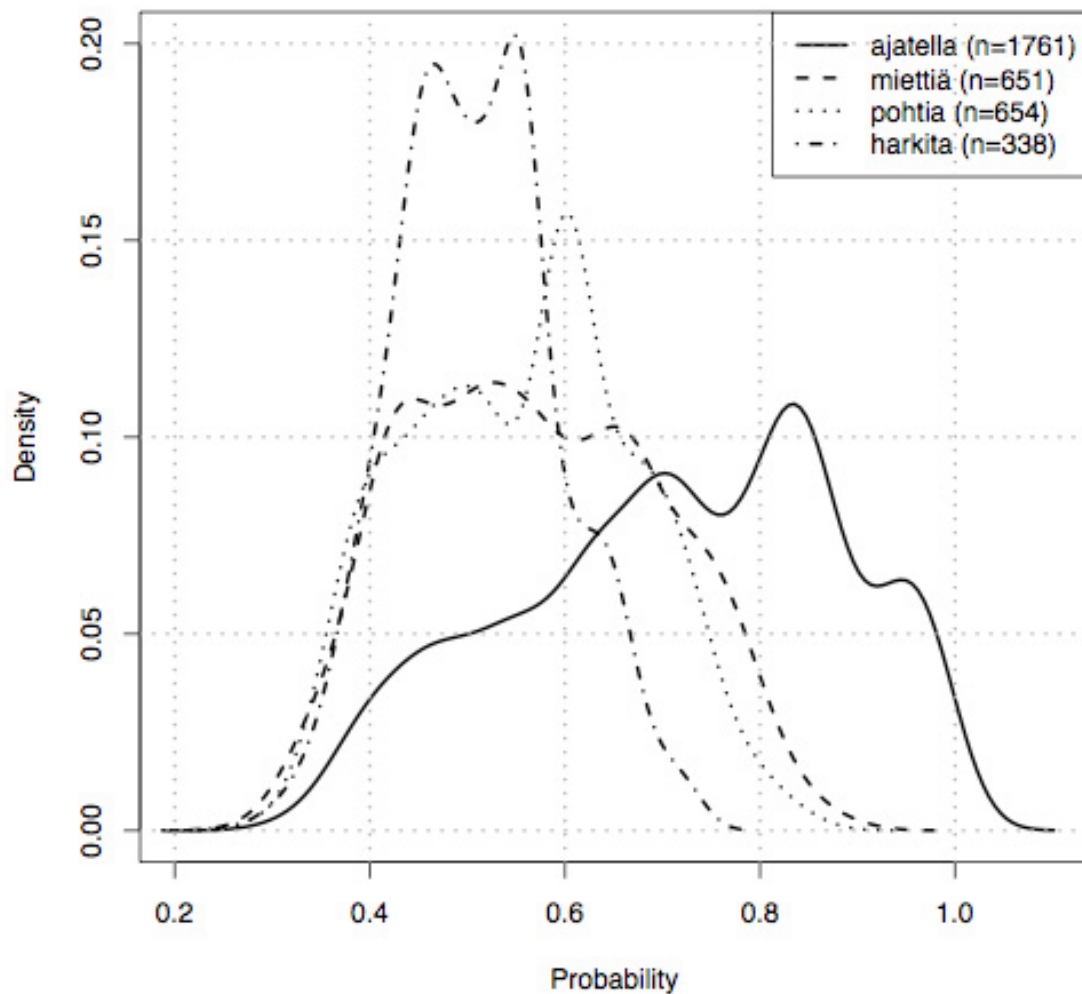


Figure 5.7. Densities of the distributions of the estimated probabilities by each lexeme according to their selection per each instance in the data.

An alternative interesting angle is what are the distributions of the probability estimates for the lexemes that originally do occur in the data set, calculated on the basis of the contextual features that are present at each such instance. In Figure 5.8 representing these particular probability distributions, we can see that the entire probability range is evident and in use. However, the bulk of the occurrences of *ajatella* receive $P > 0.5$, with the maximal peak just above $P \approx 0.8$, whereas the rarer lexemes can be assigned almost any value between zero and $P \approx 0.8$. Again, the rarer lexemes have multiple peaks, which for *miettä* are roughly at $P \approx 0.7$, $P \approx 0.4$, and $P \approx 0.2$, for *pohtia* at $P \approx 0.6$, $P \approx 0.4$, $P \approx 0.2$, and $P \approx 0.1$, and for *harkita* at $P \approx 0.5$ and $P \approx 0.1$. The maximal peaks for the rarer lexemes are not exactly in the order of their overall frequency, as for *miettä* the mode is just below $P \approx 0.4$, while for *pohtia* it is higher at just above $P \approx 0.6$, but for *harkita* at as low as $P \approx 0.1$, which entails that original occurrences of *pohtia* are predicted at a relatively high confidence, second only to *ajatella*.

Overall, these results would again suggest that some particular contexts and the associated feature combinations are relatively frequent among the lexemes, leading to the observed peaks in their estimated probabilities. Furthermore, as in this four-outcome setting any lexeme-wise estimated probability estimate $P < 0.25$ by definition amounts to its non-selection in the particular instance, the peaks below that value might indicate contexts for which the model exhibits its least accurate performance, which thus concerns roughly half of the occurrences of *harkita* and a smaller but still substantial proportion in the case of both *mieltiä* and *pohtia*. These results are concordant with the precision values for lexeme prediction presented earlier in Table 5.10 in Section 5.2.2.

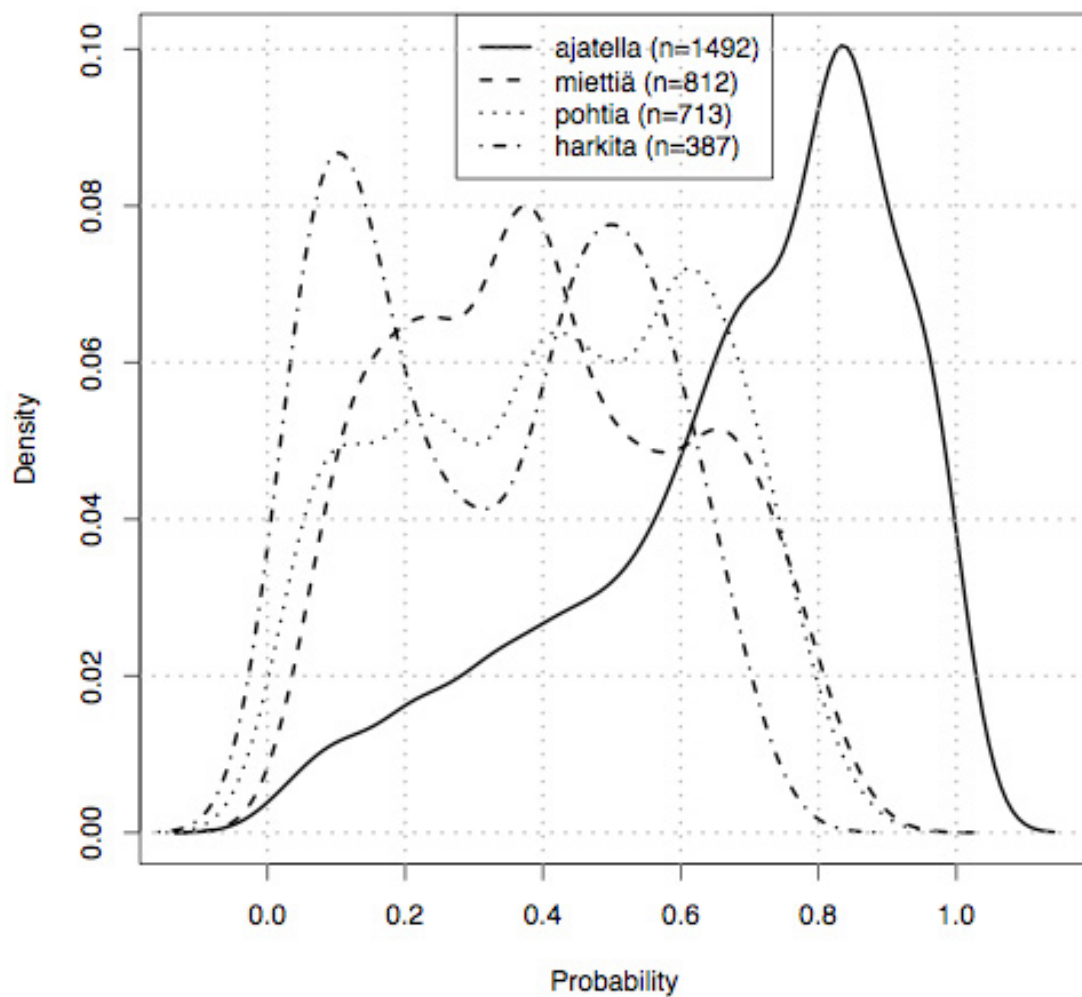


Figure 5.8. Densities of the distributions of the estimated probabilities by each lexeme according to their original occurrence per each instance in the data.

The distributions of estimated probabilities can also be represented by their partition into probability bins, which will be of assistance in the stratified selection of subsets of the original sentences and their associated feature combinations, whether as example sentences for lexicographical descriptions, or as raw materials for follow-up studies, for example, for use as experimental stimuli (see, e.g., Bresnan 2007). In Table 5.30, I

have partitioned the entire possible (adjusted) probability range $P=[0,1]$ into 10 bins with equal intervals, though in many cases five bins might practically be fully sufficient. As can be seen, the bins for *ajatella* are relatively evenly populated save for the lowest bin ($0.0 \leq P < 0.1$). However, for the rarer THINK lexemes the spread is more skewed, as among the three highest bins ($P \geq 0.7$) for *harkita* the lowest one is sparse and the two higher ones fully empty, while the situation is only slightly better for *miettiä* and *pohtia*, with the second-highest bins sparse and the highest ones altogether empty for these two lexemes.

Table 5.30. Frequencies of instances (contextual feature combinations) in the research data set ($n=3404$) for which the lexeme-wise adjusted probability estimates fall into probability bins based on 10 equal intervals.

THINK.multivariate.one_vs_rest.verb_chain_morphology.syntax_semantics_selected.extra.probability_bins 10

Probability range/Lexeme	ajatella	miettiä	pohtia	harkita
$0.0 \leq P < 0.1$	788 (23.1%)	938 (27.6%)	1453 (42.7%)	2351 (69.1%)
$0.1 \leq P < 0.2$	385 (11.3%)	906 (26.6%)	633 (18.6%)	494 (14.5%)
$0.2 \leq P < 0.3$	243 (7.1%)	497 (14.6%)	407 (12.0%)	156 (4.6%)
$0.3 \leq P < 0.4$	264 (7.8%)	400 (11.8%)	288 (8.5%)	83 (2.4%)
$0.4 \leq P < 0.5$	200 (5.9%)	225 (6.6%)	187 (5.5%)	139 (4.1%)
$0.5 \leq P < 0.6$	207 (6.1%)	172 (5.1%)	142 (4.2%)	132 (3.9%)
$0.6 \leq P < 0.7$	323 (9.5%)	147 (4.3%)	215 (6.3%)	44 (1.3%)
$0.7 \leq P < 0.8$	333 (9.8%)	104 (3.1%)	70 (2.1%)	5 (0.1%)
$0.8 \leq P < 0.9$	403 (11.8%)	15 (0.4%)	9 (0.3%)	0
$0.9 \leq P \leq 1.0$	258 (7.6%)	0	0	0

5.5.2 Profiles of instance-wise distributions of the lexeme-wise probability estimates

Finally, we can use the probability estimates to rank and select example sentences that best embody and represent the feature contexts in which the studied THINK lexemes are most typically used. Table 5.31 below contains the highest-ranked sentence for each of the lexemes which also actually contains in the original corpus data the lexeme assigned the highest probability in such a context, as well as the estimates for the other lexemes and the number of features (counting both the total and only the robust ones) which played a role in the calculation of the estimates.⁹⁴ As we can see, the sentence with the highest probability estimate for *miettiä* did *not* in fact contain this lexeme in the original data, to which issue I will return later on, so here I had to go for the second-highest ranked sentence. Examining the lexeme-specific winning sentences, we can note that the combination of GENERIC subtype of a MANNER argument, INDICATIVE mood, and SECOND person (SINGULAR) manifested in the node-verb, no explicitly expressed AGENT (i.e., COVERT) though it may implicitly be deduced on the basis of the person/number feature to be an INDIVIDUAL, and an INFINITIVE as PATIENT, together yield the maximum probability $P=1.0$ for the occurrence of *ajatella*, and accordingly also zero probabilities for the other three

⁹⁴ N.B. The source/subcorpus is counted as a feature only when the instance in question originates from the Internet newsgroup discussion, i.e., when the associated selected variable (Z_EXTRA_SRC_sfnet) is TRUE, even though its value as FALSE unequivocally determines the source as the newspaper subcorpus; the same applies also for quoted citations within the newspaper subcorpus.

THINK lexemes.⁹⁵ This can be considered as an example of a case where one or more features in the context result in a categorical choice. The overall number of features considered in the estimation thus adds up to 7, of which 2 belong to the robust set.

In turn, a CO-ORDINATED VERB (belonging to the MENTAL subtype, though this detail is not included in the multivariate regression Model VIII), FREQUENCY (representing the SOMETIMES subtype, also excluded), SECOND person (and by association also the IMPERATIVE mood) expressed by the node verb, no explicit AGENT (COVERT), though the morphology of the node-verb determines it as INDIVIDUAL, and an INDIRECT QUESTION as PATIENT jointly produce the highly preferential but not quite categorical estimate of $P=0.878$ for *miettiä* to occur, whereas among the other lexemes *pohtia* has the next highest estimate $P=0.084$, while the chances for the other two are close to zero. Thus, with these lexeme-wise estimates potential variation as well as (non-categorical) preferences/tendencies in choice of lexeme in the context are now evident, since while we could expect *miettiä* in such a context approximately nine times out of ten, as also is the case here, we could also expect *pohtia* once in every ten times, and in principle both *ajatella* and *harkita*, too, but very rarely.

In fact, this particular combination of features occurs only once in its entirety in the research corpus, which is this instance in question, so the estimates for the other lexemes are produced by weighing in occurrences of the features in other combinations throughout the research corpus. Accordingly, there are 13 instances in the data in which any five of the aforementioned six (linguistic) features are present, for which cases *miettiä* is the predominant lexeme with 11 (84.6%) occurrences, while *pohtia* is only occasional with its 2 (15.4%) occurrences. Thus, we can see that by slightly relaxing the contextual setting the outcome proportions come to roughly equal to the estimated probabilities for the entire set with six linguistic contextual features. Moreover, this example demonstrates that the instance-wise context-based probability estimates are not merely an artefact resulting from applying a probabilistic method to the data, but correspond to actual proportions of outcomes evident in the data (which logistic regression in particular aims to model).

Moving on to the maximal estimated probabilities for the last two lexemes, a (physical) LOCATION, a GROUP as an AGENT, INDICATIVE mood, THIRD person as well as PLURAL number expressed in the node verb, and a NOTION as PATIENT lead to a probability of $P=0.852$ for *pohtia* to occur, while all three other lexemes receive non-nil estimates but with clearly $P<0.1$. Finally, a clause-adverbial META-argument, an INDIVIDUAL as AGENT, CONDITIONAL mood, and NECESSITY expressed in the verb-chain, a POSITIVE (specifically THOROUGH) evaluation of MANNER, and ACTIVITY as PATIENT give *harkita* an estimated (adjusted) $P=0.725$.⁹⁶ However, neither *miettiä* nor *pohtia* are improbable in this context, as their estimates are $P(\textit{miettiä}|\textit{Context})=0.115$ and $P(\textit{pohtia}|\textit{Context})=0.135$. Thus, we have here an example of three out of the four selected THINK lexemes each having a reasonable chance of occurring in the context in question. The probability estimates for all four THINK lexemes for all the sentences in the research data, as well as their lexeme-wise rankings, the numbers of overall and

⁹⁵ One must remember, however, that the actual probabilities are not fully this categorical, as the exact values are $P(\textit{miettiä}|\textit{Context})=1.48e^{-08}$, $P(\textit{pohtia}|\textit{Context})=2.36e^{-09}$, and $P(\textit{harkita}|\textit{Context})=2.94e^{-09}$, though in practice these are, of course, as good as nil.

⁹⁶ Indeed, the corresponding unadjusted probability estimate for this particular context, from the binary model pitting *harkita* against the rest, is considerably higher at $P=0.914$.

robust features considered in the estimation, and the sentences themselves, have been compiled compactly together into one data frame `THINK.dictionary.data`, incorporated in the `amph` data set. A small selection of these containing the five highest ranked sentences per each lexeme is presented in Tables R.23-26 in Appendix R.

Table 5.31. Highest ranked example sentences (in terms of the expected probability estimates according to their contextual feature set) per each THINK lexeme, which are also matched with the occurrence of the same lexeme in the original data.

Ranking ($n_{\text{features,all}}/n_{\text{features,robust}}$) Probability estimates	Sentences
A:#1 (7/2) $P(\text{ajatella} \text{Context})=\underline{1}$ $P(\text{mieltiä} \text{Context})=0$ $P(\text{pohtia} \text{Context})=0$ $P(\text{harkita} \text{Context})=0$	<i>Miten</i> _{MANNER+GENERIC} ajattelit _{INDICATIVE+SECOND, COVERT,} <i>AGENT+INDIVIDUAL erota</i> _{PATIENT+INFINITIVE} <i>mitenkään jostain SAKn</i> <i>umpimielisistä luokka-ajattelun kannattajasta?</i> [SFNET] [3066/politiikka_9967] ‘How did you think to differ at all from some uncommunicative supporter of class-thinking in SAK?’
M:#2 (7/1) $P(\text{ajatella} \text{Context})=0.018$ $P(\text{mieltiä} \text{Context})=\underline{0.878}$ $P(\text{pohtia} \text{Context})=0.084$ $P(\text{harkita} \text{Context})=0.02$	<i>Vilkaise</i> _{CO-ORDINATED_VERB(+MENTAL)} <i>joskus</i> _{FREQUENCY(+SOMETIMES)} <i>valtuuston esityslistaa ja mielti</i> _{(IMPERATIVE+)SECOND, COVERT,} <i>AGENT+INDIVIDUAL monestako</i> _{PATIENT+INDIRECT_QUESTION} <i>asiasta sinulla</i> <i>on jotain tietoa.</i> [SFNET] [2815/politiikka_728] ‘Glance sometimes at the agenda for the council and think how many issues you have some information on.’
P:#1 (6/3) $P(\text{ajatella} \text{Context})=0.036$ $P(\text{mieltiä} \text{Context})=0.071$ $P(\text{pohtia} \text{Context})=\underline{0.852}$ $P(\text{harkita} \text{Context})=0.041$	<i>Suomessa</i> _{LOCATION(+LOCATION)} <i>kansalaisjärjestöt</i> _{AGENT+GROUP} pohtivat _{INDICATIVE+THIRD+PLURAL} <i>uudenmuotoisen auttamisen</i> <i>periaatteita</i> _{PATIENT+NOTION} (<i>mm. A-tilaajan tunnistus</i>) <i>ns.</i> <i>puhelinauttamisen eettisessä</i> <i>neuvottelukunnassa</i> _{LOCATION(+GROUP)} . [1259/hs95_10437] ‘In Finland civic organizations are pondering the principles of novel forms of assistance (e.g., the identification of an A-subscriber) in the so-called ethical advisory board of telephone assistance.’
H:#1 (7/2) $P(\text{ajatella} \text{Context})=0.025$ $P(\text{mieltiä} \text{Context})=0.115$ $P(\text{pohtia} \text{Context})=0.135$ $P(\text{harkita} \text{Context})=\underline{0.725}$	<i>Monen puoluetoverinkin mielestä</i> _{META} <i>esimerkiksi Kauko</i> <i>Juhantalon</i> _{AGENT+INDIVIDUAL} <i>olisi</i> _{CONDITIONAL+THIRD} <i>pitänyt</i> _{VERB_CHAIN+NECESSITY} harkita <i>tarkemmin</i> _{MANNER+POSITIVE(<THOROUGH)} <i>ehdokkuuttaan.</i> [275/hs95_2077] ‘In the opinion of many fellow party members, for instance, Kauko Juhantalo should have considered more carefully his candidacy.’

5.5.3 “Wrong” choices in terms of lexeme-wise estimated probabilities

Returning to the issue of the lexeme-wise probability estimates not matching the actually selected THINK lexeme in the original data, Table 5.32 exhibits the example sentence for each lexeme with the highest estimated probability which in fact contains the “wrong” lexeme. In the correspondence of the probability estimates with the actual choices the four lexemes clearly differ, as out of the 10 highest ranked sentences for both *ajatella* and *pohtia* all ten in each case contain the original lexeme, whereas the conformance level is 8/10 for *mieltiä* and 7/10 for *harkita*. Extending the

window to the 100 highest ranked sentences per each lexeme, the level of correspondence between the lexeme assigned the context-wise highest estimated probability and the original lexeme choice remains high for *ajatella* with 97/100 matches, whereas for *pohtia* this figure has dropped to 83/100, while the accuracy levels of 70/100 for *miettiä* and 62/100 for *harkita* have not substantially slipped further down.

Each of the selected example sentences in Table 5.32 present different scenarios. The first sentence contains two features with a strong preference for *ajatella*, namely, SOURCE as an argument and the AGREEMENT subtype of MANNER, resulting in a very high, almost categorical estimated probability for this lexeme ($P=0.984$), while the actually occurring *miettiä* is considered quite improbable with an estimated $P=0.014$, that is, between once or twice in every hundred similar contexts. Accordingly, my linguistic intuition would find *ajatella* fully acceptable in the sentence, if not even better than the original *miettiä*, at least in the limited context that is shown. For the second sentence, after the clear but not categorical preference of *miettiä* with $P=0.889$ resulting from DURATION as a argument and an INDIRECT QUESTION as PATIENT, there are in fact two alternative lexemes with a roughly equal likelihood of occurrence, of which interestingly *pohtia* with the slightly lower assigned probability ($P=0.043$) has been selected in the original text, instead of *harkita* ($P=0.058$). Nevertheless, all three lexemes do feel acceptable in my judgement.

The third sentence presents a primarily two-way selection, since although LOCATION as an argument, NOTION as PATIENT, and a TEMPORAL expression in the verb-chain together produce a clear preference for *pohtia* ($P=0.77$), the actually selected *miettiä* is also assigned a non-negligible likelihood ($P=0.20$), that is, once every five times, in the same context. Lastly, the fourth sentence in Table 5.32 presents a clearer three-way selection than was evident in the second sentence. Now, after the most preferred *harkita* which is assigned a fairly high $P=0.725$ in the presence of a clause-adverbial META-argument, a PATIENT as ACTIVITY, CONDITIONAL mood as well as NECESSITY in the verb chain, in addition to the POSITIVE subtype of MANNER, both *miettiä* and *pohtia* (of which the latter has actually occurred) divide the remaining probability equally, receiving each the substantial estimates of $P=0.125$, that is, once every eight times. Here, too, I find all three alternative lexemes with the higher probabilities as acceptable.

Table 5.32. Highest ranked example sentences (in terms of the estimated probabilities according to their contextual feature set) per each THINK lexeme, with another (“wrong”) lexeme selected instead in the original text; highest probabilities in **boldface**; estimate for originally occurring lexeme underlined.

Ranking($n_{\text{features,all}}/$ $n_{\text{features,robust}}$) Probability estimates	Sentences
A:#23 (6/2) $P(\text{ajatella} \text{Context})=0.984$ $P(\text{mieltiä} \text{Context})=0.014$ $P(\text{pohtia} \text{Context})=0.003$ $P(\text{harkita} \text{Context})=0$	<i>Olen_{INDICATIVE+FIRST} itse_{AGENT+INDIVIDUAL} mieltinyt hieman samansuuntaisesti_{MANNER+AGREEMENT} noista motiiveista_{SOURCE}, varsinkin jos katsoo BKT ja BKT:n kasvuprosentin (nämä saa vaikka CIA World Fact Bookista) perusteella, kuka on kuka taloudellisesti lähiaikoina.</i> '[I] have myself thought somewhat similarly about those motives, especially if one looks up the GDP and the GDP growth percent (...)' [3397/politiikka_20553]
M:#1 (10/3) $P(\text{ajatella} \text{Context})=0.01$ $P(\text{mieltiä} \text{Context})=0.889$ $P(\text{pohtia} \text{Context})=0.043$ $P(\text{harkita} \text{Context})=0.058$	<i>Jos vielä_{DURATION(+OPEN)} sorrnun_{INDICATIVE+FIRST, COVERT} joskus_{TMP+INDEFINITE} pohtimaan voisiko_{PATIENT+INDIRECT QUESTION} islamisteilla tai afrikkalaisilla olla jotain omaa tuottamusta omaan ahdinkoonsa, olen varmaan jotain aivan käsittämättömän paha ja kuvottavaa, suorastaan pahuuden akselin kannatinlaakeri? [3004/politiikka_6961]</i> 'If [I] yet succumb some time to pondering whether Islamists or Africans have some of their own doing in the plight, I am surely ...'
P:#19 (6/3) $P(\text{ajatella} \text{Context})=0.018$ $P(\text{mieltiä} \text{Context})=0.2$ $P(\text{pohtia} \text{Context})=0.77$ $P(\text{harkita} \text{Context})=0.012$	<i>Volonté_{AGENT+INDIVIDUAL} tarkkailee asioita ja ilmiöitä kuin olisi_{CONDITIONAL+THIRD} pysähtynyt_{VERB-CHAIN+TEMPORAL} Palermoon_{LOCATION(+LOCATION)} mieltimään_{INFINITIVE3} tosissaan_{MANNER+OTHER} rikoksen ja rangaistuksen ongelmaa_{PATIENT+NOTION} sivistisyhteiskunnassa. [846/hs95_8122]</i> 'Volonté observes issues and phenomena as if [he] had stopped in Palermo to ponder in earnest the problem of crime and punishment in civilized society.'
H:#2 (8/2) $P(\text{ajatella} \text{Context})=0.025$ $P(\text{mieltiä} \text{Context})=0.125$ $P(\text{pohtia} \text{Context})=0.125$ $P(\text{harkita} \text{Context})=0.725$	<i>Tarkastusviraston mielestä_{META} tätä ehdotusta_{PATIENT+ACTIVITY} olisi_{CONDITIONAL+THIRD, COVERT} syytä_{VERB_CHAIN+NECESSITY} pohtia tarkemmin_{MANNER+POSITIVE}. [766/hs95_7542]</i> 'In the opinion of the Revision Office there is reason to ponder this proposal more thoroughly.'

5.5.4 Contexts with lexeme-wise equal probability estimates – examples of synonymy?

The preceding scrutinies have already given some indication of potential interchangeability, in other words, some degree of synonymy, among one or more of the studied THINK lexemes in given contexts, though always with a clear, predominant preference for one individual lexeme. We can pursue this to the extreme and select contexts in which estimated probabilities for all four lexemes are as equal as possible. This can be measured in terms of the value range of the probabilities (i.e., the difference between the maximum and minimum probabilities per instance $\max[P(\text{Lexeme}|\text{Context})]-\min[P(\text{Lexeme}|\text{Context})]$), or their standard deviations (σ), the latter having been used in Table 5.33 to select five example sentences for which the linguistic contextual information would not appear to be able to produce

substantial distinctions between the four THINK lexemes as to their probability of occurrence.

In all but two of the sentences in Table 5.33 I could quite easily accept the substitution of the original THINK lexeme with any of the three others without any reservations. But even the remaining instance, namely, the use of *harkita* in the first and the second sentences, could be considered acceptable after creatively imagining suitable extra-linguistic circumstances, being for this particular case the implicit assumption of the thinking process to concern choices regarding actions, general behavior or opinions. We can also note that in four sentences out of five the lexeme assigned the highest probability was not actually selected in the original text, something we could naturally expect as the estimated probabilities are not that different to begin with. Consequently, these particular cases can be considered examples of contexts in which the studied THINK lexemes are as an entire set mutually most interchangeable, that is, synonymous with each other (according to a contextual definition of the concept), and overall as corroborating evidence for considering the selected THINK synonyms as near-synonyms.

Table 5.33. Example sentences for which the lexeme-wise probability estimates in terms of the contextual feature sets are most similar (on the basis of the standard deviation σ of the probabilities); highest probabilities in **boldface**; estimate for originally occurring lexeme underlined.

$n_{\text{features,all}}/n_{\text{features,robust}}(\sigma)$ Probability estimates	Sentences
5/1 (0.038) P(<i>ajatella</i> Context)=0.201 P(<i>mieltiä</i> Context)= 0.282 P(<i>pohtia</i> Context)= <u>0.279</u> P(<i>harkita</i> Context)=0.238	<i>Korkalaisella on itsellään ollut vaikea lonkkavamma ja hän_{AGENT+INDIVIDUAL} on_{INDICATIVE+THIRD} pohtinut paljon_{QUANTITY(+MUCH)} vammaisuuden kohtaamista_{PATIENT+ACTIVITY}. [3185/hs95_8865] 'Korkalainen himself has had a difficult hip injury and he has pondered a lot facing disability.'</i>
5/1(0.039) P(<i>ajatella</i> Context)= <u>0.255</u> P(<i>mieltiä</i> Context)= 0.273 P(<i>pohtia</i> Context)= 0.28 P(<i>harkita</i> Context)=0.193	<i>Suurimmat vammat saa lapsena tiukkaan lahkoon kuulunut, joka_{AGENT+INDIVIDUAL} on_{INDICATIVE+THIRD} joutunut_{VERB-CHAIN+NECESSITY} ajattelemaan_{MANNER(+LIKENESS)} lahkoon tavalla_{MANNER(+LIKENESS)} saada_{REASON/PURPOSE} rakkautta äidiltä. [2790/hs95_7550] '... who has had to think like the sect in order to receive love from [one's] mother.'</i>
8/1 (0.044) P(<i>ajatella</i> Context)= 0.301 P(<i>mieltiä</i> Context)=0.272 P(<i>pohtia</i> Context)=0.215 P(<i>harkita</i> Context)=0.212	<i>Aluksi harvemmin, mutta myöhemmin tyttö alkoi viettää öitä T:n luona ja vuoden tapailun päätteeksi P_{AGENT+INDIVIDUAL} sanoi, että voisi_{CONDITIONAL+THIRD,VERB-CHAIN+POSSIBILITY,COVERT} ajatella asiaa_{PATIENT+ABSTRACTION(<NOTION)} vakavamminkin_{MANNER+POSITIVE-(SFNET)} [50/ihmissuhteet_8319] '... P said that [he] could think about the matter more seriously [perhaps]'</i>
5/2 (0.047) P(<i>ajatella</i> Context)=0.256 P(<i>mieltiä</i> Context)=0.183 P(<i>pohtia</i> Context)= <u>0.27</u> P(<i>harkita</i> Context)= 0.291	<i>Siwan löydös on nyt tuonut_{VERB-CHAIN+EXTERNAL, VERB-CHAIN+NECESSITY} pohdittavaksi_{CLAUSE-EQUIVALENT} myös_{META} muita mahdollisuuksia_{PATIENT+ABSTRACTION(<NOTION)}. [3361/hs95_14185] 'The Siwa find has now raised for consideration also other possibilities.'</i>
4/2 (0.050) P(<i>ajatella</i> Context)=0.221 P(<i>mieltiä</i> Context)= 0.317 P(<i>pohtia</i> Context)= <u>0.259</u> P(<i>harkita</i> Context)=0.203	<i>Tuorein pohtittava_{CLAUSE-EQUIVALENT,VERB_CHAIN+NECESSITY} asia_{PATIENT+ABSTRACTION(NOTION)} on pääsihteerin ehdotus YK:n valmiusjoukkojen luomiseksi (QUOTE). [3160/hs95_2086] 'The most recent issue to be considered is the secretary-general's proposal to create a UN rapid deployment force.'</i>

Scrutinizing the actual linguistic contexts in the example sentences in Table 5.33, I find it difficult to identify any additional contextual features or essentially new feature categories, pertaining to current, conventional models of morphology, syntax, and semantics that are not yet incorporated in the current analysis at least to some extent but which would allow for distinguishing among the lexemes or selecting one above the rest, at least in the immediate sentential context. It seems rather that the semantic differences between using any of the THINK lexemes in these example sentences are embedded and manifested in the lexemes themselves, and these distinctions are of the kind that do not and would not necessarily have or require an explicit manifestation in the surrounding context and argument structure. That is, the selection of any one of the THINK lexemes in these sentences each emphasizes some possible, though slightly distinct aspect or manner of thinking, though all such aspects could be fully conceivable and acceptable as far as the constraints set by the surrounding linguistic structure are concerned. In this, the relevant discriminatory selective characteristics would concern features outside the traditional linguistic domain, that is, the expressed attitude, emotion, and style, the “nuances” which Inkpen and Hirst (2006: 1-4) have found surprisingly apt in reduplicating which of the various near-synonymous

alternative lexemes (with the tested sets comprising more than two synonyms) have actually been used, with accuracy levels even exceeding 90 percent (Inkpen and Hirst 2006: 26-27; see also Inkpen 2004: 111-112).

Take, for example, the four variations below (4.1-4.2) of the third sentence above (#3 in Table 5.33), with *ajatella* as the originally selected lexeme. When *ajatella* is used in this context (4.1), the implication to me is that the AGENT (*P*, apparently a female on the basis of the overall context, which is verified in the preceding text) might consider the PATIENT, that is, *asia* ‘matter, issue’, as a more serious affair, and therefore, change her general attitude to or disposition vis-à-vis to the matter. On the other hand, selecting *mieltiä* (4.2) conveys rather that *P* might give the matter some moments of (dedicated, if brief) thought for some unspecified duration and frequency, whereas *pohtia* (4.3) would indicate giving the matter serious, intense, and possibly lengthy consideration. Finally, if *harkita* were selected (4.4), this would mean that the matter involves some decision or choice (or abstaining from such an action) that would be reached as a result of the thinking process. Though none of these shades of meaning, which could be considered to incorporate the implications and presuppositions discussed by Hanks (1996), can be resolved on the basis of the immediate sentence context alone, they might be deduced from prior passages in the same text from which the particular sentence is taken, or from previous related texts in the same thread of discussion, or on the basis of extralinguistic knowledge about the context or even concerning the participant persons in the linguistic exchange (cf. Hanks 1996: 90, 97). Nevertheless, this case and the others with roughly equal estimates of probability represent in my view the explanatory limits of linguistic analysis which can be reached within immediate sentential context and by applying current, conventional theories and models.

- (4.1) [*Sitä sitten seurasi vuoden tapailu.*] *Aluksi harvemmin, mutta myöhemmin tyttö alkoi viettää öitä T:n luona ja vuoden tapailun päätteeksi P_{AGENT+INDIVIDUAL} sanoi, että voisi_{CONDITIONAL+THIRD,VERB-CHAIN+POSSIBILITY,COVERT} **ajatella** asiaa_{PATIENT+ABSTRACTION(<NOTION)} vakavamminkin_{MANNER+POSITIVE}.*
 ‘[That was followed by a year of dating.]⁹⁷ At first, only occasionally, but then later the girl started spending nights at T’s place and after a year of dating P said that [she] could think of the matter more seriously [perhaps]’
- (4.2) ... *P_{AGENT+INDIVIDUAL} sanoi, että voisi_{CONDITIONAL+THIRD,VERB-CHAIN+POSSIBILITY,COVERT} **mieltiä** asiaa_{PATIENT+ABSTRACTION(<NOTION)} vakavamminkin_{MANNER+POSITIVE}*
 ‘... P said that [she] could give the matter some thought [at some time or another] more seriously, [perhaps].’
- (4.3) ... *P_{AGENT+INDIVIDUAL} sanoi, että voisi_{CONDITIONAL+THIRD,VERB-CHAIN+POSSIBILITY,COVERT} **pohtia** asiaa_{PATIENT+ABSTRACTION(<NOTION)} vakavamminkin_{MANNER+POSITIVE}*
 ‘... P said that [she] could think over [at length, with concentration] the matter more seriously [maybe]’
- (4.4) ... *P_{AGENT+INDIVIDUAL} sanoi, että voisi_{CONDITIONAL+THIRD,VERB-CHAIN+POSSIBILITY,COVERT} **harkita** asiaa_{PATIENT+ABSTRACTION(<NOTION)} vakavamminkin_{MANNER+POSITIVE}*
 ‘.. P said that [she] could consider [her view with respect to] the matter [and what to do about it consequently] more seriously, [perhaps].’

⁹⁷ This preceding sentence has been added as it renders the following passage grammatical and semantically complete, as the word *harvemmin* ‘occasionally’ as well as the clause initiated by *mutta* ‘but’ in the scrutinized sentence refer back to *tapailu* ‘[occasional] dating’ in the preceding sentence, specifically its increasing frequency over time. The oddness of the selected sentence on its own was noted to me by my father Juhani Arppe.

Another recent example exhibiting such contextually non-explicit differences of meaning can be found in Figure 5.9, from the comic strip *Fingerpori* by Pertti Jarla, published daily in Helsingin Sanomat, this particular one on 8.2.2008. In this exchange between the central character Heimo Vesa and his wife Irma⁹⁸, she, astonished by her husband's non-existent table manners – he has grabbed a jelly roll and is proceeding to munch it with his bare hands, – asks Heimo *ajatteletko ikinä muita?* ‘Do you [don't you] ever think about others [other people]’. In responding with *no ... lähikaupan myyjää joskus* ‘Well ... the saleswoman at the local shop, sometimes’, Heimo follows a syntactically possible but pragmatically multiply awkward interpretation. Though the results in this dissertation have shown that human INDIVIDUALS as PATIENT arguments prefer *ajatella*, the query and its reference to the object of thought manifested typically by the PATIENT argument should in this particular extralinguistic context (i.e., the crass behavior visually evident in the strip and even more so the husband-and-wife relationship between the two characters) be understood rather as ‘Do you ever take others [other people] **into consideration**’.



Figure 5.9. A contextually non-explicit, semantically ambiguous use of *ajatella*, as exhibited in *Fingerpori* (by Pertti Jarla, © Punishment Pictures and PIB) in Helsingin Sanomat on 8.2.2008.

Among the 3404 sentences in the research corpus, there are 26 instances in which the difference between the maximum and minimum estimated probabilities $P_{max} - P_{min} \leq 0.2$, that is, all four probability estimates fall within such a narrow span. Interestingly, these sentences represent both subcorpora in exactly equal proportions (13 both), so neither text type would appear more prone to synonymous usage than the other. In these sentences, out of the altogether 46 feature variables included in the proper full model (excluding extra-linguistic ones), 23 occur at least once and 16 at least twice. The most common such features associated with structurally synonymous usage are INDIVIDUAL as AGENT ($n=16$), THIRD person ($n=14$), ABSTRACTION as PATIENT ($n=10$), INDICATIVE ($n=8$) as well as CONDITIONAL ($n=8$) mood, usage as a CLAUSE-EQUIVALENT form ($n=8$), and ACTIVITY as PATIENT ($n=8$). What is somewhat surprising is that while some of these features have been judged as neutral with respect to the lexemes in the multivariate analysis (e.g., THIRD person), most have been identified as having significant odds in favor of or against one or more of the studied lexemes (e.g., ABSTRACTIONS and ACTIVITIES as PATIENT).

⁹⁸ Identities ascertained by Pertti Jarla, the cartoonist himself (Personal communication 18.2.2008).

5.5.5 Deriving general scenarios of probability distribution profiles with clustering

In conclusion, we have been able to observe various scenarios of how the estimated probability space can be distributed among the studied THINK lexemes per individual instances on the basis of the selected features manifested in each context. Firstly, the probability distribution may approach categorical, exception-less choice, so that, in practice, only one of the lexemes is assigned the maximum possible probability $P \approx 1.0$, while the rest receive none. Secondly, selectional situations for some contexts may inherently incorporate variation so that one lexeme is clearly preferred in such circumstances, receiving the highest probability, but one or more of the others may also have a real though occasional chance of occurring to a varying degree. This was shown to logically result in individual instances of actual usage for which the selected lexeme is not the one which was assigned the highest probability estimate. Lastly, we have also observed cases in which all four lexemes are estimated to have approximately equal probability with respect to the observable context, as it can be linguistically analyzed according to current, conventional theory, so that any differences in meaning are conveyed by the particular selected THINK lexeme alone.

In addition to these somewhat accidentally identified though quite sensible *impromptu* scenarios, however, we can in fact apply statistical clustering techniques to systematically arrange and group the entire set of lexeme-wise probability distribution estimates now available for all instances ($n=3404$) in the data. Using hierarchical agglomerative clustering (HAC), with the *Euclidean* distance measure and the *Ward* clustering algorithm, and next selecting with simplicity first in mind a quite arbitrary quantity of five clusters from the result,⁹⁹ we can then calculate, for each cluster, mean values for the probability estimates over whichever lexeme happens to receive the maximum, second-highest, third-highest, and minimum frequency at a time.¹⁰⁰

As can be seen in Figure 5.10, cluster 2 corresponds most to the practically categorical scenario, with one lexeme receiving almost all of the available probability (with an average $P \approx 0.89$), leaving very little to the other three (all $P \leq 0.07$). Clusters 1, 2, and 5 can be considered exemplars of the variable but preferential outcome case, where more than outcome is in practice possible and to be expected, occasionally. However, whereas cluster 1 represents a two-way choice, where one lexeme is clearly preferred over the second one (with $P \approx 0.71$ vs. $P \approx 0.17$), with the other two relegated as marginal, cluster 5 exhibits a three-way outcome scenario, where the two alternatives rated second-highest and third-highest are relatively equal (with $P \approx 0.24$ and $P \approx 0.17$, respectively), though the highest-rated lexeme still stands above them taking the majority with $P \approx 0.53$. Furthermore, cluster 2 presents a variation of the former in which the lexeme rated second-highest comes considerably closer to the most highly rated one (with $P \approx 0.57$ vs. $P \approx 0.34$), these two becoming in practice the only viable alternatives in comparison to the remaining other two ($P \leq 0.06$). Finally,

⁹⁹ Other numbers of clusters could potentially be more motivated on the basis of a more thorough scrutiny and analysis of the results of the hierarchical clustering algorithm, but as this is not the main focus in this dissertation I have decided to opt for a preliminary, tentative treatment of the question with an exploratory character, which may be refined in later research.

¹⁰⁰ Thus, in this clustering process no distinction is made concerning which of the individual THINK lexemes receive the highest, lowest, or any other rankings of the probability estimates for each instance; rather, the focus is on the general instance-wise distribution of probability estimates.

cluster 4 comes closest to the synonymous case but not quite; though none of the lexeme-wise probabilities receives overall predominance with $P > 0.5$, their range remains nevertheless quite broad, being equally spaced between $P \approx (0.09, 0.042)$. More detailed expositions of the ranges (minima, maxima as well as the 25% and 75% quartiles in addition to the means) of the probability estimate values falling under each cluster, presented in Figures R.1-5 in Appendix R, affirm that these clusters represent distinct probability estimate profiles. Thus, with the help of a statistical technique we have been able to both verify and generalize the prior instance-specific analyses as well as bring forth new details concerning the studied phenomenon.

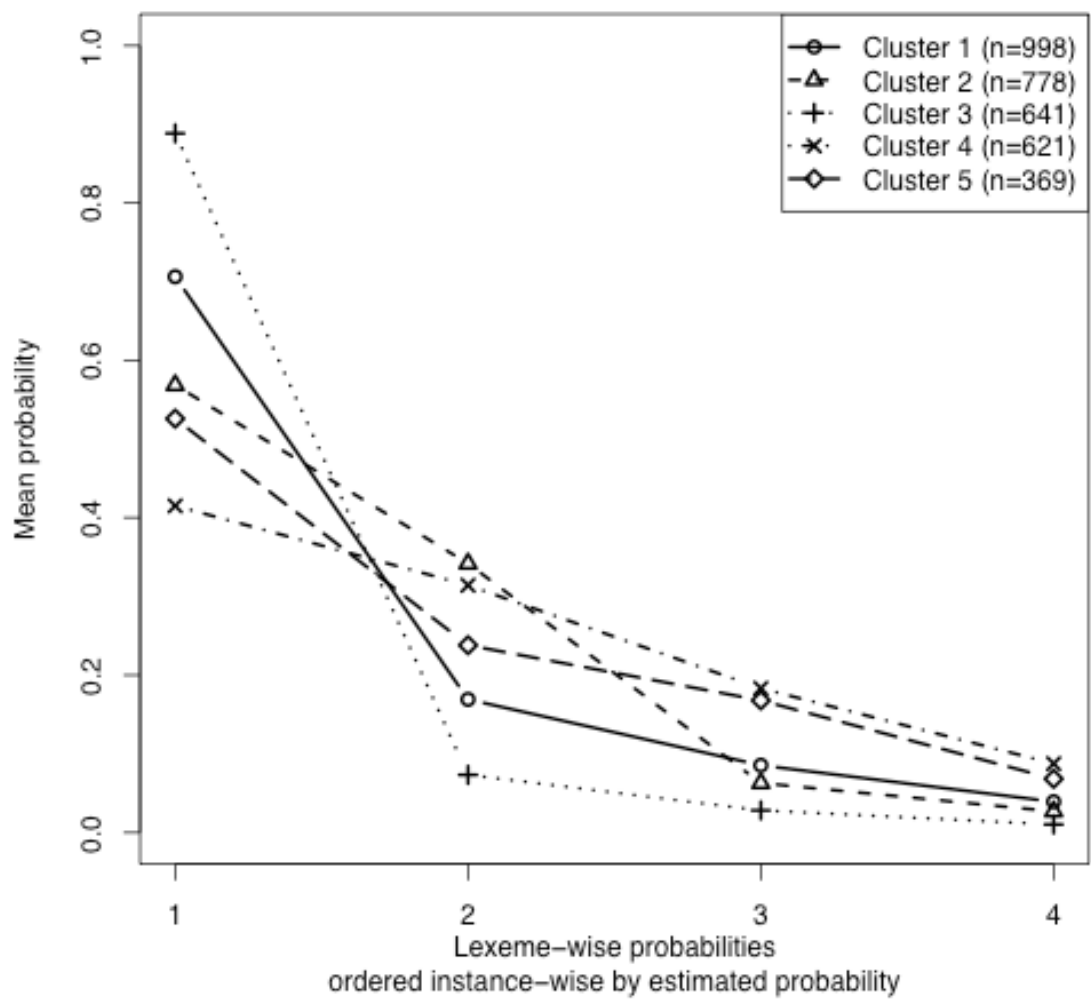


Figure 5.10. Lexeme-wise mean probabilities, in descending order, for five clusters of instance-wise distributions of probability estimates in the research data set.

5.6 New descriptions of the studied synonyms

We can close off this extensive discussion of the results with some thoughts concerning how they could in practice be used to draw up new descriptions of the studied synonyms, such that could be used in dictionary entries in works such as *Perussanakirja* (PS). The corpus-based identification of contextual features which distinguish the synonyms from each other and the feature-wise lexeme-specific odds which were assigned to them with the polytomous logistic regression analysis, especially those in favor of a lexeme occurring, would quite naturally form the basis of a more formal description of the synonyms and their usage. The features in favor of a lexeme could in the first place be listed, ordered according their individual descending odds. The sentences in the research corpus on which the analyses are based, coupled with the lexeme-wise probability estimates that are the second essential output of the polytomous logistic regression analysis, can readily be seen as the raw material from which exemplary, real usage contexts can be selected.

The major obstacle to using these current results directly and picking the sentences with the highest probability estimates for each lexeme is that similar contexts will receive similar probability estimates; thus, such straightforward selection will lead to examples which are essentially duplicates of the same contexts and features. We can correctly assume that for each lexeme there are several typical usage contexts, which cannot necessarily be reduced to one individual genuinely observed sentence that would aggregate them all (cf. Divjak and Gries' [2006: 42] similar judgement concerning the improbability of singular examples incorporating all the criteria associated with Idealized Cognitive Models), even though there probably will be substantial overlap among them as the lexemes individually and as a whole do share some common contextual features, for example, prototypically HUMAN beings as AGENT. Consequently, a good practical description would contain exemplars for each of such distinct contexts, for which one could also hope to be able to incorporate several relevant (robust) features at a time for the sake of economizing as well as adhering to reality. Furthermore, if one indicates the contextual features in such selected example sentences, one can also render the formal description more intelligible to non-professional users of such a description.

Since the entire set of sentences in the research corpus has already been classified according to the contextual feature variables which we have used throughout the above analysis, we can use the sentence-wise feature sets as input for a statistical clustering algorithm to sort the underlying sentences into groups which are internally similar but group-wise distinct. Hierarchical agglomerative clustering (HAC), which has already been applied in this dissertation, is an attractive clustering technique as it does not require us to determine the number of clusters beforehand; in contrast, it allows us to extract whatever number of clusters we may deem practical later on. Since the sentence-wise features are binary, the corresponding *Binary* distance measure seems most appropriate. With respect to the various clustering algorithms available, the *Single linkage* method, which adopts a "friends of friends" clustering strategy, ends up clustering practically all the sentences together into a single, all-encompassing group, which can be considered an indication of the relative overall semantic and contextual overlap among the sentences. Consequently, the *Ward* clustering algorithm, which aims at finding compact, spherical clusters, is in principle more advantageous, and it was found to produce useful clusters of more or less

similar size. On the basis of this clustering, having first determined an overall number of clusters (either arbitrarily or on the basis of more closer, additional scrutiny of the clustering structure), we can next pick from each cluster one or more examples according to a range of criteria, such as the estimated probability of the lexeme occurring in the sentence, the overall number of features, or the number of robust features present in the context manifested in the sentence.

```
THINK.data.verb_chain_morphology.syntax_semantics_semantics.extra.clu
stered_binary_single
THINK.data.verb_chain_morphology.syntax_semantics_semantics.extra.clu
stered_binary_ward
```

Setting the number of clusters at 50, their sizes range between 13 and 190 sentences, averaging some 68 sentences. Though these clusters to some extent appear to be associated with one or two individual lexemes, with 13 clusters for which the most frequent lexeme accounts for over two-thirds of the sentences, nevertheless 39 clusters contain at least one exemplar of all 4 lexemes and 9 such exemplars of any three. Picking the probability-wise most highly-ranked example per cluster turns out to be strongly biased towards the most frequent of the set, *ajatella*, coming on top 45 times out of 50. Using a more elaborate selection scheme emphasizing “richness” with respect to features, when we pick one sentence for each of the 50 clusters which have been sorted, first, by their number of robust features and, second, by their overall number of features, and third, in terms of estimated probability, we can extract a set of examples consisting of 20 sentences with *ajatella*, 16 with *miettiä*, and 13 with *pohtia*, but only 1 with *harkita*. Variations of this selection strategy can be tried out with the `select.sentences.by.clusters` function.

If one would like to further offset the general preference of *ajatella* in terms of the probability estimates and instead compile a comprehensive set of examples for each lexeme which would at the same time be mutually maximally distinct, an alternative strategy would be to pick from all the clusters for each lexeme the sentence assigned the highest probability, with the requirement that the lexeme in question has a genuine, substantial chance of occurring in the contexts represented by each cluster (i.e., $P \geq 0.5$). Keeping the number of clusters at 50, this lexeme-oriented procedure produces a total of 113 sentences, of which 43 are with *ajatella*, 28 with *miettiä*, 23 with *pohtia*, and 19 with *harkita*, which is more in line with the overall proportions of these lexemes in the research data. This latter selection strategy can be tried out and varied using the `select.sentences.by.lexemes_and_clusters` function. The two sets of example sentences according to the different aforementioned strategies, along with the lexeme-wise probabilities, are stored in the data tables `THINK.dictionar.y.selection.robust_feature_probability_10.50` and `THINK.dictionar.y.selection.lexemes_by_clusters.p.min_.5.50`, respectively, in the `amph` data set.

In conclusion, a new description, for example, for *pohtia*, following the latter mentioned selection strategy, with morphological, syntactic and semantic preferences and a representative set (23) of example sentences, is presented in Table 5.34. Similar descriptions for the other studied THINK lexemes can easily be compiled along the same lines. We can see that all but one of the preferred features are present in abundance among the examples, with the exception of expressions or media of COMMUNICATION as PATIENT, a specific case which could be added to the example set

by hand. Overall, this result can in fact be considered a close variant of the *Behavioral profiles* as presented by Hanks (1996). While the example sentences with the explicit indication of the relevant contextual features can be considered to combine Hanks' formal complementation patterns with genuine usage, the underlying clustering as well as emphasis of the number of features alongside high expected probability ensures that all essential arguments and their combinations are represented, that is, the "totality of their complementation patterns" (Hanks 1996: 77-78), though their number is greater than the (in general) maximally dozen found sufficient by Hanks (1996: 84).

Furthermore, the expected probability of a sentence and the features it incorporates has in Table 5.34 replaced the relative frequency used in the behavioral profiles by Hanks (e.g., 1996: 80, Figure 2) as an indication of typicality. Nevertheless, building directly on the classification of the original contextual elements, the examples in Table 5.34 lack deeper interpretations of intentions and presuppositions of the type of hypothesized characterizations presented in Table 4.5 in Section 4.1.2, which also Hanks (1996: 90, 97) concedes may not be extractable efficiently by computational means, but rather on the basis of (possibly collective) introspection by linguists or lexicographers. In the end, whereas Hanks (1996: 84-85) characterizes his behavioral profiles as the building blocks of a "dictionary without definitions", I would describe the description presented in Table 4.34 as a stepping stone towards a *dictionary of examples*.

Table 5.34. A new description of *pohtia* with 1) a formal presentation of its contextual preferences as well as 2) a set of representative example sentences, with the preferred features (i.e., ones with significant odds >1 in the multivariate analyses) indicated with subscripts.

Features	PATIENT+DIRECT_QUOTE (8.1) AGENT+ GROUP (4.2) PATIENT+ABSTRACTION (4.1) LOCATION (3.7) PATIENT+COMMUNICATION (3) PATIENT+INDIRECT_QUESTION (2.8) VERB-CHAIN+TEMPORAL (2.4) TIME-POSITION+DEFINITE (2.3) PASSIVE (1.9) PATIENT+ACTIVITY (1.6) PLURAL (1.6)
Examples	(0.852) Suomessa _{LOCATION} kansalaisjärjestöt _{AGENT+GROUP} pohtivat _{PLURAL} uudenmuotoisen auttamisen periaatteita _{PATIENT+ABSTRACTION} (mm. A-tilaajan tunnistus) ns. puhelinauttamisen eettisessä neuvottelukunnassa. [hs95_10437] (0.844) Pari lehteä _{AGENT+GROUP} ehti _{VERB-CHAIN+TEMPORAL} jo sunnuntaina _{TIME-POSITION+DEFINITE} pohtimaan pääkirjoituspalstoillaan _{LOCATION} valtion vakuusrahaston johtajan Heikki Koiviston ennenaikaista eroamista _{PATIENT+ACTIVITY} . [hs95_2140] (0.815) Hän neuvoi viimeaikaisiin tapahtumiin viitaten, että EU:ssa _{LOCATION} ryhdyttäisiin _{VERB-CHAIN+TEMPORAL} pohtimaan keinoja _{PATIENT+ABSTRACTION} rajoittaa "siirtolaisuutta islamilaisista maista". [hs95_2786] (0.811) ... lis. Osmo Soininvaara pohtivat _{PLURAL} yksilön ja yhteisön sosiaalista vastuuta _{PATIENT+ABSTRACTION} klo 19 _{TIME-POSITION+DEFINITE} ravintola Kahdessa Kanassa _{LOCATION} , Kanavakatu 3, Katajanokka. [hs95_9522] (0.806) Lohjan kunnan sosiaalidemokraattien, oikeiston ja keskiryhmien, vihreiden ja vasemmistoliiton valtuustoryhmät _{AGENT+GROUP} pohtivat _{PLURAL} kuntaliitosasioita viikonvaihteessa _{TIME+POSITION+DEFINITE} . [hs95_9607]

(0.800) Asiaa_{PATIENT+ABSTRACTION} **pohdittiin**_{PASSIVE} viime viikolla_{TIME-POSITION+DEFINITE} Helsingissä UNHCR:n järjestämässä suljetussa seminaarissa_{LOCATION}. [hs95_10142]

(0.782) Tarvetta_{PATIENT+ABSTRACTION} muuttaa vahingonkorvauksia aletaan_{VERB-CHAIN+TEMPORAL} **pohtia** oikeusministeriön asettamassa työryhmässä_{LOCATION}. [hs95_2890]

(0.756) Hän pitää käytännössä mahdollomana, että maailman kaikki YK:n sopimuksen solmineet valtiot_{AGENT+GROUP} saataisiin_{PASSIVE} koolle **pohtimaan** tapaus_(PATIENT+EVENT) Estoniaa. [hs95_7496]

(0.733) Kuvassa Juha Kankkunen (takana) **pohtimassa** rengasvalintaa_{PATIENT+ACTIVITY} RAC-rallissa_{LOCATION} 1992_{TIME-POSITION+DEFINITE}. [hs95_4892]

(0.732) Iltapäivällä_{TIME-POSITION+DEFINITE} **pohditaan**_{PASSIVE} ryhmissä_{LOCATION} kehitysysteistyötä_{PATIENT+ACTIVITY}, liikennettä, maataloutta, suomalaista luontoa, ympäristöä ja kulutusta sekä energiaa. [hs95_1154]

(0.723) Viimeksi suomalaiset teatterintekijät **pohtivat**_{PLURAL} noin runsas puoli vuotta sitten_{TIME-POSITION+DEFINITE} Tampereen teatterikesässä_{LOCATION}, miksi_{PATIENT+INDIRECT_QUESTION} varsinkin monet naisohjaajat haluavat tarkastella elämän ikuisia peruskysymyksiä juuri myyttien näkökulmasta. [hs95_10041]

(0.714) Suomen kulttuurista tulevaisuutta_{PATIENT+ABSTRACTION} **pohtimaan** tänään nimitettävä Maanantaiseura_{AGENT+GROUP} on saamassa jäsenikseen paitsi poliittisia konkareita myös uuden sukupolven nimiä: listalle on kaavailtu niin puoluejohtaja Ulf Sundqvistia kuin City-lehden toimittajaa Eero Hyvöstä. [hs95_8772]

(0.688) “Tuntuu siltä, että lännen arkkitehtien varakkuus on heidän suurimpia vaarojaan, sillä mukavuudet katkaisevat yhteyden ikuisen luontoon”_{PATIENT+DIRECT_QUOTE}, virolainen arkkitehti Leonhard Lapin **pohtii** näyttelynsä saatesanoissa_{LOCATION}. [hs95_9762]

(0.683) Lasten vieminen ja hakeminen päiväkodista on useimmiten isän kontolla, ja päiväkodin tädit ovat tahollaan_{LOCATION} **pohtineet**_{PLURAL}, miten_{PATIENT+INDIRECT_QUESTION} kommunikoida tämän lasta kuljettavan miehenkilon kanssa. [hs95_15267]

(0.657) Yleisönosastossa_{LOCATION} on alettu_{VERB-CHAIN+TEMPORAL} **pohtia**, mistä_{PATIENT+INDIRECT_QUESTION} puulajista saa parhaan leipälapion. [hs95_71]

(0.654) Tässäkin näyttelyssä_{LOCATION} on oiva tilaisuus **pohtia** suomalaisen modernismin myöhäsyntyistä olemusta_{PATIENT+ABSTRACTION}, josta C. J. af Forselles kirjoittaa pitkään luettelossa. [hs95_1794]

(0.651) Vähimmäisturvan ja muiden sosiaalirahojen tason yhtenäistäminen ei kuulu toimikunnan_{AGENT+GROUP} **pohdittaviin** asioihin_{PATIENT+ABSTRACTION}. [hs95_15473]

(0.643) “... Me täällä humanistisessa tiedekunnassa_{LOCATION} **pohditaan**_{PASSIVE} oikeasti tärkeitä asioita_{PATIENT+ABSTRACTION} (kuten maitolaitureiden epävirallinen käyttö vuosina 1950- 60) niinku todella syvällisesti ja silleen...” [ihmissuhteet_1060]

(0.640) EDUSKUNNAN perustuslakivaliokunta_{AGENT+GROUP} on tehnyt tarkkaa työtä **pohtiessaan** ministerien jääviyden rajoja_{PATIENT+ABSTRACTION}. [hs95_2143]

(0.552) Niemi ja Kotikumpu **pohtivat**_{PLURAL} odotusaitiossa_{LOCATION} suksien vaihtoakin_{PATIENT+ACTIVITY}, mutta se ei käynyt pänsä erilaisten siteiden takia. [hs95_12402]

(0.532) Ei kai näitä sisältöjä_{PATIENT+ABSTRACTION} voi **pohtia** tanssin jytkeessä_{LOCATION}. [hs95_10082]

(0.514) Ei siksi, että miehetkin **pohtisivat**_{PLURAL} nyt suhteita_{PATIENT+ABSTRACTION}, vaan siksi, että maskuliinisuuksien kenttä jää kokoelmassa hahmottomaksi ja hajanaiseksi. [hs95_9757]

(0.514) Rautiainen **pohti** asioita_{PATIENT+ABSTRACTION} harvinaisen kokonaisvaltaisesti ja asetti “riman” korkealle myös itselleen. [hs95_960]

This description for *pohtia* in Table 4.34 can now be compared with the ones currently available in *Perussanakirja* (PS) and *Nykysuomen sanakirja* (NS) presented earlier in Tables 2.9 and 2.10 in Section 2.3.2. In comparison to both PS and NS, with respect to the form of exposition this new description makes the preferred contextual features explicit by having the associated lexemes or structures marked in the example sentences. If one were to filter out all but the specifically preferred feature context, the results would, to some degree, resemble the truncated model phrases presented in both PS and NS, but they are lengthier and contain typically more than one feature at a time, and thus, perhaps, rather similar to the citations provided in NS.

In terms of linguistic content, it is interesting to note that the THOROUGH (POSITIVE evaluation) subtype of MANNER presented in both PS and NS is not to be found among the corpus-derived examples, as if this connotation, which is real for me as a native speaker of Finnish, were incorporated in the lexeme (*pohtia*) itself, not requiring an overt exponent in the context. In fact, this context is more particular to another THINK lexeme, namely, *harkita*. Similarly, CO-ORDINATION as well as QUANTITY (specifically its MUCH subtype) both apparent in NS are attributed instead to *mieltiä* according to the corpus-based results. Furthermore, PURPOSE or REASON as syntactic arguments, exemplified in association with *pohtia* in both PS and NS, are judged to exhibit only a dispreference in conjunction with *ajatella* on the basis of the research corpus, which also links the expression of NEGATION in the verb-chain with *ajatella*, with a significant dispreference for *pohtia*, thus contrasting one example in NS. However, ABSTRACTIONS and ACTIVITIES as PATIENT both occur in Table 4.34 as well as in PS and NS; likewise, an INDIRECT QUESTION as PATIENT, LOCATION as a syntactic argument (though its ABSTRACTION subtype presented in NS is rare in the research corpus, with $n=11$), and the PASSIVE voice are apparent both in Table 4.34 and NS.

In contrast, the description in Table 4.34 contains a range of features altogether absent from PS, and also to a slightly lesser degree from NS, namely, GROUPS as AGENT, DIRECT QUOTES as PATIENT, the DEFINITE subtype of TIME-POSITION, and the expression of TEMPORALITY (specifically, START) and PLURAL number in the verb-chain. Consequently, the differences between the new description provided in Table 4.34 and the earlier ones presented in both PS and NS are substantial, but neither are the corpus-derived results reached in this dissertation entirely discordant with the contents of the earlier dictionaries. Nonetheless, I do hope that the results encapsulated in Table 4.34 will contribute to fulfilling the duty which Atkins and Levins (1995: 107) place upon linguists – and linguistics as a discipline – to provide the theoretical infrastructure on which lexicographical descriptions can be soundly based and improved.

6 Discussion

6.1 Synonymy and its study and description in light of the results

On the whole, the results of the analyses presented in this dissertation demonstrate the great variety of the different feature categories and the complexity of their interrelationships that must be grasped in order to explain the studied synonym group of Finnish THINK verbs. Yet, it must be said that the results also indicate the limits of conventional linguistic analytical features in this endeavor. First of all, the univariate results (Section 4.1) show that a wide range of different linguistic (morphological, lexical, syntactic and semantic) and extralinguistic (text type, medium and repetition) contextual features are individually associated with the usage and context-wise appropriate selection of the chosen four-member synonym set, which is in line with the observations and conclusions of Divjak and Gries (2006). Many contextual features were observed in the research corpus to have conspicuous and clearly significant preferences or dispreferences for the studied lexemes which were not exemplified at all in the respective word entries in the latest authoritative dictionaries of Finnish, *Suomen kielen perussanakirja*, that is, PS (Haarala et al. 1997) or *Nykysuomen sanakirja*, that is, NS (Sadeniemi 1976 [1951-1961]), nor in the more formal overall description of the Finnish verb system by Pajunen (2001). In addition, a substantial number of features were incorporated among the examples in PS or its predecessor, NS, which were analyzed to have a dispreference with respect to the particular lexemes in question on the basis of their distribution in the research corpus, thus indicating a further discrepancy between the current lexicographical descriptions and actual usage as exhibited by the research corpus. Secondly, the bivariate results (Section 4.2) indicate that the features are pairwise interconnected to a varying but for the most part quite weak degree.

Thirdly, the multivariate results (Section 5), based on polytomous logistic regression modeling, show that taken together the features have different weights and importance in determining which of the lexemes, and with what anticipated probabilities of occurrence, can be expected to be used in a particular context incorporating a given set of features. By and large, there is more than one solitary feature, identified as statistically relevant with respect to the studied lexemes, extant in their intrasentential contexts; in fact, the median number of such contextual features per instance in the research corpus is 5, with a lower quartile (25%) of 4 and an upper quartile (75%) of 6, and a maximum of as many as 11 features. Nevertheless, a few features which may appear individually significant in univariate analyses can turn out *not* to play a significantly distinctive role in multiple feature considerations. Furthermore, though morphological features, either concerning just the node verbs or the entire verb-chain of which they form a part, exhibit clear preferential distinctions among the studied THINK lexemes in the univariate analyses (as was already observed in Arppe 2002 and later in Arppe and Järvikivi 2007b), their overall importance is in the end diminished in the multivariate analysis in comparison to the semantic classifications of syntactic arguments or the verb-chain as a whole, the latter two feature categories which receive the highest weights (confirming the similar initial observation in Arppe 2007). A similar fate with an even stronger drop in relative importance applies to extralinguistic variables, when they are considered together with the other feature categories in the multivariate analysis. Moreover, the most stringent assessment of the robustness of the results, bootstrapping with resampling from speakers/writers as

clusters, indicated that only about one-half of the observed preferences showed pervasiveness over the entire population represented in the selected data, thus suggesting the potential for generalization beyond the studied research corpus.

Viewed from the standpoint of the estimated probabilities for lexical outcomes given a set of contextual features, the results indicate that there exists for the most part substantial and tangible variation with respect to which lexemes can actually occur in the close-to-same contexts (Section 5.5). In fact, for 77.6% of the sentences in the research corpus the estimated expected probabilities are for all four lexemes at least $P(\text{Lexeme}|\text{Context}) > 0.01$. This variation can be categorized into several general scenarios, of which the most characteristic ones are noted here. Firstly, the observed proportions and the estimated probability distributions may approach categorical, exception-less choice, so that only one of the lexemes is assigned in practice the maximum possible probability, while the rest have nil probability. However, such a scenario applies to as few as 7.6% of the sentences in the research corpus. Secondly, many contexts may inherently incorporate variation so that one of the lexemes is clearly preferred in such circumstances, receiving by far the highest probability, but one or more of the others may also have an occasional but nevertheless tangible chance of occurring. Sometimes, two lexemes may account among themselves for almost all of the observed occurrences in some particular context, with the other two in practice not occurring at all. Lastly, we have also observed cases in which all four lexemes are estimated to have approximately equal probability vis-à-vis the observable context and linguistic features applied in this dissertation.

These instances with close-to-equal estimated probabilities of occurrences could be considered prime candidates as examples of “genuine” synonymy and complete interchangeability in context (Section 5.5.4). Scrutinizing the linguistic contexts of such sentences, I found it difficult to identify any additional contextual features or essentially new feature categories which would allow for distinguishing among the lexemes or selecting one over the rest, at least in the immediate sentential context. Rather, it seems that the semantic differences between using any of the THINK lexemes in these example sentences are embedded and manifested in the lexemes themselves, and these distinctions are of the kind that do not and would not necessarily have or require an explicit manifestation in the surrounding context and argument structure. That is, the selection of any one of the THINK lexemes in these sentences each emphasizes a possible – though slightly distinct – aspect or manner of thinking, which are all contextually equally acceptable and fully conceivable. Nevertheless, these distinctions could possibly be deducible from the entire text or chain of associated texts, or even the overall extralinguistic context, if such contexts were available to an observer.

Overall, the *Recall* rate of correctly predicting lexical choice among the four selected THINK lexemes seems to reach a ceiling at approximately two-thirds, or 64.6-65.6% to be exact, of the instances in the research corpus (Section 5.2.3), and appears to be indifferent to whether an individual group of variables is left out or the variable set is substantially increased (ignoring for the moment the recommended limitations to the size of the feature set with respect to the minimum frequencies of outcomes in the data). The question that again first springs to mind is whether we still lack some necessary variables or variable types, perhaps pertaining to discourse or information structure, which have been applied in prior studies with dichotomous selectional

settings (e.g., Gries 2003a, 2003b, Bresnan 2007). Or, one might suspect that the more complex polytomous setting scrutinized in this study is more difficult to accurately model, though the different heuristics used to implement this have for all practical purposes produced equal results in terms of prediction performance. Moreover, it is also conceivable that taking the interaction effects among the feature variables – which were left out of the models in this study – into account might improve model fit and accuracy.

The closer inspection of not only sentences with roughly equal estimates of probability for all four lexemes but also those with non-categorical preferences for one or two of the lexemes suggests that such selectional variation in context is both common and acceptable, and that any distinctive features there may be are not explicitly evident in the immediate sentential context, but rather pertain to stylistic attitudes and intended shades of expression that the speaker/writer wishes to convey (pertaining to the intermediate stylistic/subconceptual level in the clustered model of lexical choice by Edmonds and Hirst 2002). Furthermore, similar, less than perfect levels of prediction accuracy (54%¹⁰¹), have been reached for the even more complex 6-way prediction of synonymous Russian TRY verbs, using the simultaneously fit multinomial heuristic with a baseline category, on the basis of the semantic properties of their subjects and the following infinitives as well as Tense-Aspect-Mood (TAM) marking on the TRY verbs themselves (personal communications from Dagmar Divjak 4.12.2007, 16.5.2008, and 19.5.2008), suggesting that the performance levels reached in this dissertation are not at all exceptionally poor or low.

In conclusion, the observed general upper limit of *Recall* in prediction, as well as the sentences with roughly equal estimates of probability can be viewed to represent the explanatory limits of linguistic analysis attainable within the immediate sentential context and by applying the conventional descriptive and analytical apparatus based on currently available linguistic theories and models (cf. Gries 2003b: 13-16). Moreover, the results also indicate that contextual (i.e., distributional) similarity would not appear to lead us to full (absolute) synonymy, that is, (full) equality in meaning. More generally, these results support Bresnan's (2007) probabilistic notion about the relationship between linguistic usage and the underlying linguistic system (see also Bod et al. 2003). Few choices are categorical, given the known context (feature cluster) that can be analytically grasped and identified. Rather, most contexts exhibit various degrees of variation as to their outcomes, resulting in proportionate choices in the long run. Since these context-relative proportions of outcomes, which logistic regression in specific aims to replicate as probability estimates, *are* evident in the data (albeit roughly), their probabilistic character cannot be dismissed as merely an artefact resulting from applying a probabilistic method to the data. Nevertheless, this probabilistic view of language is neither fully accepted yet, nor necessarily irreconcilable with the categorical view (see, e.g., Yang 2008 and references therein).¹⁰²

¹⁰¹ In the validation of this model, the jack-knife estimate was 50,8%, while randomly splitting 100-fold the entire data sample of 1351 instances into training sets of 1000 instances and testing sets with the remaining 351 instances yielded a mean correct classification rate of 49%, with a standard deviation of 2.45% (Personal communication from Dagmar Divjak 16.5.2008).

¹⁰² I am thankful to my external reviewer Stefan Th. Gries for drawing my attention to the controversy concerning this question.

From the overall methodological perspective, the three different levels of analysis, namely, the univariate, bivariate, and multivariate ones with their respective statistical methods, could each be observed to play in turn an essential role in discovering the most important explanatory features, thus supporting Gries' (2003a) general multivariate (i.e., "multifactorial" in his terminology) approach, which also entails proceeding through all these stages, starting with the simplest univariate scrutinies, followed by pairwise comparisons, and only then finishing with the most complex multivariate analyses. The goal in the univariate analysis is to identify a comprehensive range of distinct linguistic perspectives (i.e., feature categories as well as individual features) which are relevant with respect to the studied phenomenon. At this univariate level, the well-established chi-squared (χ^2) test of the homogeneity/heterogeneity of the distribution of a feature among the studied lexemes, followed up by standardized Pearson residuals (e_{ij}) for the scrutiny of individual lexeme-specific preferences for each feature, appeared to be the most useful method, already quite reliably anticipating the directions, though not the strengths, of the preferences/dispreferences to be uncovered in the later multivariate analyses (Sections 3.2.2 and 4.1.1). Furthermore, considering the distributions of multiple related features at the same time among the studied lexemes, referred to as grouped-feature analysis in this dissertation, produced for the most part similar preference/dispreference patterns in relation to singular-feature scrutinies, which compare an individual feature's occurrences against its nonoccurrences among the studied lexemes (Sections 3.2.3 and Appendix N).

Among the various summary measures of association, I found Theil's asymmetric Uncertainty Coefficient ($U_{B|A}$), belonging to the Proportionate Reduction in Error, or alternatively, Proportion of Variance Explained (PRE) category of methods, as the most useful in assessing how much individual features accounted for the variation among the studied lexemes. Still, it is important to keep in mind that the features are multiply intercorrelated in real linguistic usage so that they also, in any individual instance, are bound to incorporate the influence of the other relevant features concurrently present in the context (Sections 3.2.2 and 4.1.1). Moreover, as an asymmetric measure, $U_{B|A}$ could be used to assess to what extent a lexeme-feature preference/dispreference relationship could be regarded as either feature-specific or lexeme-specific (following Arppe and Järvikivi 2007b: 148), though in the single-feature scrutinies the lexemes were overall (per feature) always the more dominant determinant of the two possible directions in the feature-lexemes relationship. However, most of the individual features were found to account for only a small proportion of the variation in the studied phenomenon, although a few more influential features were also observed. Moreover, the selected chi-squared-based symmetric measure of association, namely, Cramér's V , was not found to have a meaningful correlation with PRE measures such as the $U_{B|A}$, underlining the non-PRE character of the former measure. Taking an additional perspective within univariate analysis (with the results presented in Appendix K), I also explored scrutinizing the distributions of the features among the studied THINK lexemes from the Zipfian perspective, but as the number of items in the selected synonym set was quite low (being only four), no really significant results were to be gained.

Turning to the bivariate analysis (Sections 3.3 and 4.2), I chose here, too, to use the Uncertainty Coefficient (denoted this time as $U_{2|1}$) in the pairwise comparisons of the distributions of features, since it is asymmetric in its values for the 2x2 setting

crosstabulating the occurrences and nonoccurrences of two features against each other, in comparison to other possible measures of association with a similar conceptual basis and interpretation, namely, the Goodman-Kruskal $\tau_{2|1}$. Furthermore, the grouped-feature scrutiny of the homogeneity/heterogeneity of distributions already applied in the univariate analysis could likewise be extended to compare two sets of individually related features, specifically in order to identify individual pairings of features which exhibit higher than expected co-occurrences.

With respect to multivariate statistical methods, firstly, polytomous logistic regression and secondly, the one-vs-rest technique for its implementation were both shown to be attractive methods in the study of lexical choice with multiple alternative outcomes, thus building upon and extending Bresnan et al.'s (2007) work which was restricted to a dichotomous alternation (Sections 3.4 and 5). More specifically, as has been noted earlier in this dissertation, logistic regression provides naturally interpretable analysis results in assigning odds for the explanatory features by which their relative importance in describing the observed phenomenon can be assessed and compared. Furthermore, logistic regression can integrate the joint occurrence of multiple contextual features, the kind often evident within a sentence in normal language usage, as one single statistic, estimating the expected probability of occurrence of an outcome in such a context, which should correspond to the originally observed proportions of outcomes in the same contexts. Furthermore, the one-vs-rest heuristic was shown to perform equally well in comparison to other, allegedly more sophisticated or "elegant" techniques, supporting Rifkin and Klautau's (2004) emphatic arguments in favor of its use. In particular, I found the one-vs-rest technique the most appealing among the various alternatives, since it provides lexeme-specific estimates of feature-wise odds for all outcome classes, that is, lexemes in this study, without the need for assuming or selecting some prototypical baseline category, not to mention its obvious practical simplicity.

In comparison to the Hierarchical Agglomerative Clustering (HCA) employed by Divjak and Gries (2006), polytomous logistic regression has the advantage of working on instance-wise combinations of features and individual outcomes rather than the overall proportions of features aggregated for each outcome class, thus, in principle, facilitating better consideration of the features' actual interactions (which can also be explicitly scrutinized, though that was not undertaken in this study). In assessing the robustness of the results, the bootstrap with resampling from speakers/writers as clusters is the procedure best adapted to identifying those features for which lexeme-wise preferences are pervasive throughout the entire selected population, and thus the strongest candidates for generalizations. However, the medium/source of the linguistic data simply as an additional extralinguistic variable, without interactions with the other features, was not observed to have a substantial impact on the preference patterns of the other linguistic features proper.

In an exploratory study such as is presented in this dissertation, the number of features evident in the data can turn out to be quite daunting. Nevertheless, the sets of features associated with the individual lexemes could be interpreted in a *post hoc* analysis to form coherent, meaningful groups, the characterizations of which transcend the individual features (Section 4.1.2). Consequently, the contextual associations of *ajatella* could as a whole be viewed as embodying temporal continuity, individuality in agency, and objects (PATIENTS), in addition to denoting the intentional state as a

subtype of THINKING, while for *miettiinä* its core semantic character appears to be temporally more definite, in addition to being personal and individual in agency. In contrast to *ajatella* and *miettiinä*, *pohtia* can be characterized as collective, impersonal in agency, and non-concrete with respect to its objects/PATIENTS (in a somewhat surprising contradiction with its concrete origins), whereas *harkita* can be linked to THINKING of an action as an object/PATIENT, which is temporally situated in the future. These characterizations could well be considered to represent the intermediate stylistic/subconceptual level in Edmonds and Hirst's (2002) clustered model of lexical knowledge. Furthermore, the general characterizations, as well as the features' preferences/dispreferences with which they are associated, could also be interpreted to incorporate and perpetuate historical vestiges of the concrete origins of the now quite abstract set of studied THINK lexemes.

Finally, we can use the multivariate results as a basis for actual lexicographical description (Section 5.6). For formal purposes, we can present the contextual features which have been identified to distinguish the synonyms from each other together with the feature-wise lexeme-specific odds which were assigned to them with the polytomous logistic regression analysis. For more informal purposes, we can exploit the lexeme-wise probability estimates, which are the second essential output of the polytomous logistic regression analysis, to select from the research corpus complete example sentences, which would be a convenient and effective way of embodying both a natural and a typical set of features for a lexeme in real usage. Using hierarchical agglomerative clustering (HAC), we can sort these sentences into groups which are internally similar but group-wise distinct, on the basis of the underlying sentence-wise feature sets. We can then pick from each cluster one or more examples according to a range of criteria, such as the estimated probability of the lexeme occurring in the sentence, the overall number of features, or the number of robust features present in the context manifested in the sentence. The resultant set of example sentences supplemented with the explicit indication of the relevant contextual features can be considered to modify Hanks' (1996: 77-78) notion of Behavioral Profile, consisting originally only of a formalized, abstracted description of the "totality of their [words'] complementation patterns", to rather be represented in the form of authentic, natural usage, in which the expected probability of a sentence and the feature set it incorporates works as an indicator of typicality. In my view, descriptions extracted and compiled in this manner can be regarded as stepping stones towards a *dictionary of examples*.

Yet, in the end, we may still be faced by a couple of nagging questions: Has this dissertation simply made explicit what a professional lexicographer can normally achieve – possible even surpass – by manually scrutinizing a (sufficient) set of (randomly sampled) concordances? Have I only made explicit the best practices which skilled lexicographers learn to follow in their work? Moreover, studies within computational linguistics concerning a similar task of word-sense classification suggest that results approaching a quality on par with that observed in this dissertation might be achieved through combinations of several levels of *automatic* linguistic analysis already available for many languages (cf. Lindén 2004). Consequently, I am convinced that it would be worthwhile to conduct a comparative follow-up study applying the methods presented in this dissertation simply on the raw output of the FI-FDG parser. This would allow us to quantitatively assess whether, and to what degree,

we might derive sufficiently useful lexical descriptions on the basis of automatic linguistic analysis alone, without the need for costly manual annotation.

6.2 Hypotheses for experimentation on the basis of the results

The now derived corpus-based results, specifically the multivariate ones, provide a solid basis for comparisons with other sorts of linguistic evidence, for example, experimentation such as forced-choice and acceptability rating tasks, extending and fine-tuning to multiple outcomes the simple dichotomous setting presented in Arppe and Järvi­kivi (2007b: 152, Table 5, see also Section 1.2 in this dissertation). Particularly useful and in fact quite central in these cross-evidential comparisons will be the key characteristic of polytomous logistic regression modeling which allows for the aggregation of the occurrences of varying sets of multiple contextual features into one single statistic, namely, an estimate of expected probability that should approximate the observed proportions in these same contexts in the original research data. As the vast majority of the studied features are distinctive either in favor of or against the occurrence of the individual studied lexemes, and only very few of these features are overall neutral, it would be difficult to construct experimental “laboratory” sentences for which all but one experimental variable would be controlled and neutralized, so that such sentences would also have at least some resemblance to real language usage. Single-argument, or even two-argument sentence fragments with no other words would hardly appear genuine to experimental subjects, and, what is more, too obviously underliningly indicative of the object of research. In contrast, the estimated expected probabilities allow us to take into account simultaneously a multiple of possibly occurring variables, a setting which corresponds considerably better with the makeup and composition of sentences and utterances which naturally used, produced, and encountered.

Consequently, the first and most straightforward assumption and hypothesis concerning relationships between different evidence types is that such corpus-based expected probability estimates for the entire set of lexemes in some contexts should be matched by similar proportions of the selections of these lexemes in the same contexts in forced choice tasks, given now four alternative choices instead of the two in Arppe and Järvi­kivi (2007b). This would follow from the general conclusion suggested in Arppe and Järvi­kivi (2007b) that forced-choice tasks would largely correspond to the production of corpus content as a linguistic process. So, in a practically categorical case such as sentence #1 in Table 6.1 below, with altogether seven contextual features of which two are robust, we would expect the alternative to *ajatella* to be selected almost always and the three other lexemes seldom if at all, as $P(\textit{ajatella}|\textit{Context}_{\#1})\approx 1$. In a case exhibiting variation but a clear preference for one of the lexemes such as sentence #2 in Table 6.1, with altogether seven features of which only one is robust, we would expect to see *mieltä* selected roughly nine times out of ten and *pohtia* once every ten times, with only sporadic selections of either *ajatella* or *harkita*. In contrast, in the case of structural synonymy such as sentence #3 in Table 6.1, with as many as eight features present of which only one is robust, our assumption would be to observe each of the four alternatives selected roughly equally often.

Table 6.1. A small selection of sentences from the research corpus with varying distributions of estimated probabilities for the four studied THINK lexemes, based on the results presented in Tables 5.31-5.33 in Sections 5.5.2-5.5.3; highest probabilities in **boldface**; estimate for originally occurring lexeme underlined.

#/ (Features) Probability estimates	Sentence
#1 (7/2) P(<i>ajatella</i> Context)= <u>1</u> P(<i>mieltiä</i> Context)=0 P(<i>pohtia</i> Context)=0 P(<i>harkita</i> Context)=0	<i>Miten</i> _{MANNER+GENERIC} <i>ajattelit</i> _{INDICATIVE+SECOND, COVERT,} <i>AGENT+INDIVIDUAL erota</i> _{PATIENT+INFINITIVE} <i>mitenkään jostain SAKn</i> <i>umpimielisistä luokka-ajattelun kannattajasta?</i> [3066/politiikka_9967] 'How did you think to differ at all from some uncommunicative supporter of class-thinking in SAK?'
#2 (7/1) P(<i>ajatella</i> Context)=0.018 P(<i>mieltiä</i> Context)= 0.878 P(<i>pohtia</i> Context)=0.084 P(<i>harkita</i> Context)=0.02	<i>Vilkaise</i> _{CO-ORDINATED_VERB(+MENTAL)} <i>joskus</i> _{FREQUENCY(+SOMETIMES)} <i>valtuuston esityslistaa ja mielti</i> _{IMPERATIVE+SECOND, COVERT,} <i>AGENT+INDIVIDUAL monestako</i> _{PATIENT+INDIRECT_QUESTION} <i>asiasta sinulla</i> <i>on jotain tietoa.</i> [2815/politiikka_728] 'Glance sometimes at the agenda for the council and think how many issues you have some information on.'
#3 (8/1) P(<i>ajatella</i> Context)= 0.301 P(<i>mieltiä</i> Context)=0.272 P(<i>pohtia</i> Context)=0.215 P(<i>harkita</i> Context)=0.212	[<i>Aluksi harvemmin, mutta myöhemmin tyttö alkoi viettää öitä T:n luona ja vuoden tapailun päätteeksi</i>] <i>P</i> _{AGENT+INDIVIDUAL} <i>sanoi, että</i> <i>voisi</i> _{CONDITIONAL+THIRD, VERB-CHAIN+POSSIBILITY, COVERT} <i>ajatella</i> <i>asiaa</i> _{PATIENT+ABSTRACTION(<NOTION)} <i>vakavammin</i> _{MANNER+POSITIVE-} (SFNET) [50/ihmissuhteet_8319] '[...] P said that [he] could think about the matter more seriously [perhaps]'

With respect to acceptability rating judgements, according to Bresnan (2007), we would be led to assume that such ratings would also roughly equal the estimated probabilities (when normalized to the range $P=[0,1]$), just like the proportions of selection in a forced-choice task. Firstly, however, Bresnan's results concerned a dichotomous alternation, which would not be the case with the entire set of studied THINK lexemes. Secondly, Bresnan's experiments were set up so that the binary ratings had to add up to a constant, an assumption which Arppe and Järvikivi (2007b) criticize, if Bresnan's results are taken to reflect acceptability. Therefore, I would be inclined to hypothesize on the basis of Featherston's (2005, see also Figure 1.1 in Section 1.1) results that acceptability ratings in a polytomous setting would turn out to be arranged along a linear, and only slightly descending slope, so that the highest rating would go to the lexeme assigned the highest probability, and so forth. My underlying assumption here is that due to the synonymous relationship among the studied THINK lexemes none would be considered in practice altogether unacceptable and non-interchangeable in any of the possible contexts, thus leading the lowest-judged lexemes in such contexts to nevertheless receive ratings at least in the middle range on the available scale.

In practice, it might be more manageable to consider only three of the four THINK lexemes at a time for the two forms of experimental tasks. As either *ajatella* or *harkita*, with their intentional or future-oriented uses, respectively, can each be considered the odd man out, two possible such subsets would be {*ajatella*, *mieltiä*, *pohtia*} and {*mieltiä*, *pohtia*, *harkita*}. Another way of simplifying the experimental setups would be to select contexts for which the estimated probability distributions more or less follow a specific pattern. Such patterns could be among the ones I

manually identified in Sections 5.5.2–5.5.4, and which were verified and fine-tuned with cluster analysis, for example, categorical choice with $\exists P(\text{Lexeme}|\text{Context})\approx 1$, or genuine variation but with a clear preference with $\exists P(\text{Lexeme}_1|\text{Context})\approx [0.8..0.9]$ and $\exists P(\text{Lexeme}_2|\text{Context})\approx [0.1..0.2]$, or approximate structural synonymy with $\forall P(\text{Lexeme}|\text{Context})\approx 0.25$.

6.3 Suggestions for other further research and analyses

The present study has already made considerable headway in satisfying the need for further research laid out in Arppe and Järviö (2007b: 149), namely, the extension from a synonym pair to a synonymous word group with more than two members, as well as expanding the set of contextual features considered from a few person-number morphological features and one associated syntactic argument type to the entire syntactic argument structure of the studied lexemes. However, many avenues for further research, in addition to the experimentation already discussed in Section 5.2, still remain uncharted, each which would contribute to establishing the validity and generalizability of the already achieved results.

I will first address linguistic follow-up research questions, which generally concern the generalizability of the attained results over the lexicon in individual languages as well as cross-linguistically. Firstly, similar to Divjak (2006) one could pick another related synonym group within the COGNITION verbs, such as the UNDERSTAND verbs, the most common of which in Finnish are *ymmärtää*, *käsittää*, *tajuta* and *oivaltaa* ‘understand, comprehend, grasp’, and explore what the results would be in their case, even more so as I have already studied the morphological preferences of this verb set using the visual correspondence analysis method (Arppe 2005a). One would seek to discover to what extent the syntactic argument types and their semantic classifications observed to be distinctive for the THINK verbs, for example, concerning AGENTS and PATIENTS but also others, would also figure into the context of the UNDERSTAND verbs. Likewise, one would be curious to find out which of the syntactic arguments and semantic classifications would turn up as particular to and distinctive among the UNDERSTAND verbs. Other interesting synonym sets within the COGNITION verbs could be either the TRY or INTEND verbs (or both), studied in Russian by Divjak and Gries (2006) and Divjak (2006), which would at the same time also provide an opportunity for cross-linguistic comparison. Furthermore, in this vein, one could expand the focus from the individual synonym sets to the broader semantic grouping that they belong to, namely, the COGNITION verbs, and scrutinize which features exhibit common behavior and which are distinctive for the component synonym sets within this semantic class, using perhaps only the most frequent lexeme or the aggregate of the lexemes for each synonym set. One could in a similar fashion scrutinize even the most general semantic grouping of MENTAL verbs. However, picking some synonym sets among the non-mental ACTION verbs, for example, the Finnish SHAKE/QUAKE verbs *hytistä*, *järistä*, *tutista*, *täristä*, *vapista*, *vavahdella*, *väreillä*, *värehtiä*, *väristä*, and *värähdellä*, or VERBAL COMMUNICATION verbs at the intersection of MENTAL and ACTION verbs, for example, *puhua*, *sanoa*, *kertoa* ‘speak, say, tell’, would probably most convincingly validate that the synonym group-internal distinctions occur for all types of (Finnish) verbs.

Secondly, similar syntactic-semantic contextual behavior has been observed not only within particular word classes such as verbs but also within entire morphological families derived from the same root (Argamann and Pearlmutter 2002); what is more, members of such word families have been shown to be cognitively interconnected (for an overview, see De Jong 2002). Therefore, the most common direct nominal (noun and adjective) derivations of the THINK group, for example, *ajatus* ‘thought’, *ajattelu*, ‘thinking’, *ajattelematon* ‘unthoughtful’, *miete* ‘thought’, *mietintä* ‘thinking’, *pohdinta* ‘pondering’, *harkinta* ‘consideration’, should also be investigated. For the UNDERSTAND group, similar nominal derivations would be, for example, *ymmärrys* ‘comprehension’, *ymmärtävä* ‘understanding’ (adjective), *ymmärtämätön* ‘uncomprehending’, *käsitys* ‘conception, impression’, *käsittämätön* ‘incomprehensible’, *tajuaminen* ‘realization’, and *oivallus* ‘insight’.

Thirdly, I have hypothesized in Section 4.1.2 that the roots of the current contextual preferences for the selected THINK lexemes would lie in the original concrete usages underlying the current abstract meanings. One could verify this with a study using historical corpus data concerning those of the studied lexemes which are known to have had an active concrete usage relatively recently, that is, within a few generations back, a period of time for which sufficient Finnish corpus resources are now available. This would concern at least *pohtia* and possibly *harkita*, as both words are known to have commonly used senses pertaining to farming and fishing activities, noted in Finnish dictionaries as late as the mid-1900s. Thus, the *Historical Newspaper Library*

<http://digi.lib.helsinki.fi/sanomalehti/secure/main.html?language=en> recently released by the Finnish National Library, containing scanned data from all Finnish newspapers published between 1770 and 1890, would be an attractive resource for such a historical study. Alternatively, related THINK lexemes still retaining both a concrete as well as an abstract meaning, such as *punnita* ‘weigh’, *hautoa* ‘brood, mull over, incubate’, or *märehtiä* ‘chew [over], ruminate’, which may be currently undergoing a shift towards a more abstract meaning but have yet to complete the change, could also be studied in this respect using the considerably larger and more diverse corpus resources of contemporary Finnish.

Fourthly, cross-linguistic analysis would make it possible to study whether the observed contextual preferences apply generally to languages, that is to say, are not particular only to Finnish, and to what extent this would be the case. Consequently, one could study the equivalents of the THINK lexemes in structurally divergent languages, for example, English as the “opposite” of Finnish with minimal morphology and fixed word order, and, e.g., Swedish or Russian which lie somewhere in between. A tentative hypothesis here would be that one could identify a core set of common contextual elements over languages and cultures, shared by human societies in general, while there would also be language-specific contextual preferences which could be considered culture-specific “residue”. In the case of Finnish, such culture-specific elements might again concern the rural/agricultural roots from which the Finnish society has emerged only within the latter half of the twentieth century. Furthermore, one could envision that it would be to a certain degree possible to “predict” (after the fact) contextual preferences within a synonym set on the basis of the etymologically established, original concrete usages of the individual words in

such a set.¹⁰³ Last among the linguistic extensions of this work, one could also consider applying the analysis methods presented in this dissertation to polysemy as well as synonymy, as has previously been done within the Behavioral Profile approach (Gries and Divjak, forthcoming; e.g., the study of the multiple senses of the English *run* by Gries 2006). A prime candidate for such research in Finnish would be the highly polysemous verb *pitää*, with the four main senses of ‘hold onto/hold up/keep/retain’, ‘organize/arrange’, ‘like/love’, and the modal ‘must/ought’.

As for the statistical methods presented in this dissertation, one can clearly identify two areas which would benefit from further development. In the first place, since logistic regression modeling as a multivariate method sets limits on the number of individual explanatory variables vis-à-vis the minimum desirable number of outcomes in the research data, exploring the practical implementation and the extent of improvements in the performance of statistical methods through which the set of variables could be clustered or otherwise aggregated and thus be kept at an acceptable and manageable level without losing explanatory power would be worthwhile. As has already been noted in Section 3.4.2, possible methods which should resolve this issue would be Principle Components Analysis (PCA) and related techniques, as well as Cluster Analysis, working on the intercorrelations of the explanatory variables by themselves, irrespective of the associated outcomes. A considerably simpler, alternative approach, which would be interesting to try out and test in terms of its performance, would be the selection of features for multivariate analysis lexeme-specifically for each respective individual binary logistic regression model of which the entire polytomous model is composed. However, this latter approach would not provide a solution when the number of possible features grows excessively high. Moreover, though mixed effects modeling does not, by itself, address the overabundance of variables, it is an attractive methodological development as it would allow for incorporating straightforwardly as a part of the actual statistical model longitudinal effects concerning, for instance, speaker/writer or text-specific bias.

Secondly, since I explicitly deemed visual statistical techniques beyond the scope of this dissertation, follow-up work could potentially use such methods as Correspondence Analysis or Self-Organized Maps not only to validate the now achieved results but also to provide new perspectives into the research data. Though for the current relatively low number (four) of selected THINK lexemes Hierarchical Cluster Analysis might not necessarily yield dramatically new insights as to the mutual relationships of the lexemes, it generally remains a powerful and useful tool, as was demonstrated in the case of clustering the entire set of COGNITION verbs solely on the basis of their single-word definitions and their pairwise overlaps in Section 2.1.3, even more so as its requirements with respect to minimum outcome frequencies are not as stringent as is the case for logistic regression. Furthermore, it would also be interesting to apply to the data various methods typically rather associated to the computational side of linguistics or computer science in general, such as Memory-Based Learning (MBL), Rule Induction, Random Forests, and other machine-learning and data mining techniques.

¹⁰³ I owe the concrete formulation of this line of research as well as the central hypotheses to discussions with Martti Vainio and Juhani Järvikivi.

If one wanted to build upon and generalize this work, should one first have at one's disposal a sufficiently broad general semantic ontology, of the WordNet type covering the common, core lexical content of a language arranged into synonym sets, and second, a relatively richly annotated corpus large enough to contain a sufficient number of at least the more common lexemes in the ontology, one could envision generating in an assembly-line fashion both formalized feature descriptions and representative example sentences concerning the usage of one synonym group after another, which professional lexicographers could then refine further into actual dictionary content.

7 Conclusions

In this dissertation, I present an overall methodological framework for studying linguistic alternations with multiple outcomes, focusing specifically on lexical variation in denoting a single meaning, that is, synonymy. As a practical example, I employ the synonymous set of the four most common Finnish verbs denoting THINK, namely, *ajatella*, *miettiä*, *pohtia*, and *harkita*.

Building on previous research, I describe in considerable detail the extension of statistical methods from dichotomous linguistic settings (e.g., Gries 2003a, Bresnan et al. 2007) to polytomous ones, concerning more than two possible alternative outcomes. The applied statistical methods are arranged into a succession of stages with increasing complexity, proceeding from univariate via bivariate to multivariate techniques in the end, following the general scheme laid down by Gries (2003a) in his study of a dichotomous structural alternation in English. Together, the three types of methods provide a rich overview of the phenomenon under investigation. The univariate methods can be used to identify significantly distinctive individual features with respect to the studied phenomenon, the bivariate methods can evaluate the degree of pairwise association for such features, and the multivariate methods can assess the weights and importance of individual features in relation to the entire set included in the closer examination. As the central multivariate method, I argue for the use of polytomous logistic regression and demonstrate its practical implementation to the studied phenomenon, thus extending the work by Bresnan et al. (2007), who applied simple (binary) logistic regression to a likewise dichotomous structural alternation in English. My motivation for this methodological choice is that the two main results of logistic regression modeling have natural interpretations, that is to say, in 1) the odds that are assigned for each feature incorporated in the model with respect to an outcome, indicating the increase or decrease in the chances of such an outcome occurring in conjunction with the feature in question, and in 2) the expected probabilities which can be estimated for any combination of features included in the model, approximating the actually observed proportions of outcomes in the corresponding original contexts and associated feature sets. Among the various techniques for implementing polytomous logistic regression, I find the one-vs-rest technique (Rifkin and Klautau 2004) to have the most advantages, due to its practical simplicity and descriptive characteristics, while attaining a similar performance level as other more complex and sophisticated procedures. Specifically, the one-vs-rest technique can provide in a straightforward manner feature-wise odds for *all* outcome classes – without the need for selecting a baseline, prototypical class.

As for the set of explanatory variables, I wholeheartedly agree with Gries (2003a), Divjak and Gries (2006), and Bresnan (2007) et al. (2007) in that the scientifically satisfactory and valid description of a linguistic phenomenon requires the consideration of a comprehensive range of different, relevant feature categories, and an assessment of their interactions, instead of resorting to monocausal explanations. Thus, in my analysis I incorporate feature types identified as significant and distinctive with respect to the usage and choice of synonyms in a wide range of earlier work, including lexical context (e.g., Church et al. 1991), syntactic structure (e.g., Biber et al. 1998), semantic subclasses of syntactic argument types (e.g., Atkins and Levin 1995), morphological features (e.g., Jantunen 2001, 2004; Arppe 2002), as well as text type and register (e.g., Biber et al. 1998), which corresponds to the Behavioral

Profile approach (Hanks 1996, Divjak and Gries 2006). In the linguistic analysis of the selected synonym set and their context in the research corpus, I begin with general-purpose analysis tools and resources, such as the implementation of Functional Dependency Grammar for Finnish (Tapanainen and Järvinen 1997), that is, the FI-FDG parser, on the morphological and syntactic levels, and the ontology of the English WordNet (Miller et al. 1990) for the semantic classification of (nominal) syntactic arguments, but for some less common argument types I apply the *ad hoc* evidence-driven strategy advocated by Hanks (1996). The results of the various statistical analyses confirm that a wide range of contextual features across different categories are indeed associated with the use and selection of the selected THINK lexemes; however, a substantial part of these features are not exemplified in current Finnish lexicographical descriptions. The multivariate analysis results indicate that the semantic classifications of syntactic argument types are on average the most distinctive feature category, followed by overall semantic characterizations of the verb chains, and then syntactic argument types alone, with morphological features pertaining to the verb chain and extralinguistic features relegated to the last position.

In terms of the overall performance of the multivariate analysis and modeling, the prediction accuracy seems to reach a ceiling at a *Recall* rate of roughly two-thirds of the sentences in the research corpus. Furthermore, this performance appears indifferent to whether some individual groups of feature variables are left out. Moreover, for an overwhelming majority of the sentences and associated contextual features in the research corpus, the polytomous logistic regression model in fact provides distributions of lexeme-wise estimates in which more than one lexeme is allotted genuine, tangible chances of occurring, varying from a clear but not categorical preference of one lexeme to practically equal probabilities for all four lexemes. Manually scrutinizing the linguistic contexts in various sentences in the research corpus, I found it difficult to identify any additional contextual features or essentially new feature categories, which would allow for distinguishing among the lexemes or selecting one over the others, at least within the immediate sentential context. Rather, my conclusion is that in these particular sentences the semantic differences between using any of the THINK lexemes are incorporated into and manifested in the lexemes themselves. Moreover, these distinctions are such that neither need be nor (possibly) can be expressed in some overt, explicit way in the immediately surrounding context and argument structure, even though one might – having read the entire text or knowing the overall extralinguistic context – possibly deduce these intended shades of meaning. In other words, the choice of any one of the THINK lexemes in these sentences each highlights some potential and conceivable – though slightly distinct – aspect or manner of thinking, all of which are equally acceptable with respect to the particular context.

Taken together, these last-mentioned results support Bresnan's (2007) and probabilistic view of the relationship between linguistic usage and the underlying linguistic system, in which only a minority of linguistic choices are categorical, on the basis of contextual criteria which can be observed, and thus also analyzed (understood together as a feature cluster) (see also Bod et al. 2003). Instead, most contexts exhibit degrees of variation as to their outcomes, resulting in proportionate choices over longer stretches of usage in texts or speech. Thus, the observed sentences with unequal but broadly dispersed, or even roughly equal estimates of probability represent the explanatory limits of morphological, syntactic, and semantic linguistic

analysis which we can reach within an immediate sentential context and by applying current, conventional theories and models.

Corpora

amph 2008. A micro-corpus of 3404 occurrences of the four most common Finnish THINK lexemes, *ajatella*, *mieltä*, *pohtia*, and *harkita*, in Finnish newspaper and Internet newsgroup discussion texts, containing extracts and linguistic analysis of the relevant context in the original corpus data, scripts for processing this data, R functions for its statistical analysis, as well as a comprehensive set of ensuing results as R data tables. Compiled and analyzed by Antti Arppe. Available on-line at URL: <http://www.csc.fi/english/research/software/amph/>

Finnish Text Collection [FTC] 2001. ~180 million words of Finnish, consisting of 97 subcollections of Finnish newspaper, magazine and literature texts from the 1990s. Compiled by the Research Institute for the Languages in Finland, the Department of General Linguistics of the University of Helsinki, and the Foreign Languages Department of the University of Joensuu. Available on-line at URL: <http://www.csc.fi/kielipankki/>

Helsingin Sanomat 1995. ~22 million words of Finnish newspaper articles published in Helsingin Sanomat during January–December 1995. Compiled by the Research Institute for the Languages of Finland [KOTUS] and CSC – IT Center for Science, Finland. Available on-line at URL: <http://www.csc.fi/kielipankki/>

Keskisuomalainen 1994. ~2 million words of Finnish newspaper articles published in Keskisuomalainen during January–April 1994. Compiled by the Research Institute for the Languages of Finland [KOTUS] and CSC – IT Center for Science, Finland. Available on-line at URL: <http://www.csc.fi/kielipankki/>

Parole 1998. ~16 million words of Finnish newspaper articles. Compiled by the Department of General Linguistics, University of Helsinki, and the Research Institute for the Languages of Finland [KOTUS]. Available on-line at URL: <http://www.csc.fi/kielipankki/>

SFNET 2002–2003. ~100 million words of Finnish internet newsgroup discussion posted during October 2002–April 2003. Compiled by Tuuli Tuominen and Panu Kalliokoski, Computing Centre, University of Helsinki, and Antti Arppe, Department of General Linguistics, University of Helsinki, and CSC – IT Center for Science, Finland. Available on-line at URL: <http://www.csc.fi/kielipankki/>

References

- Agresti, Alan 2002. *Categorical Data Analysis* (Second edition). Hoboken: John Wiley & Sons, Hoboken.
- Allen, James F. 1983. Maintaining Knowledge about Temporal Intervals. *Communications of the ACM*, Vol. 26, No. 11 (November 1983), pp. 832-843.
- Alonge A., N. Calzolari, P. Vossen, L. Bloksma, I. Castellon, M. A. Marti, and W. Peters 1998. The Linguistic Design of the EuroWordNet Database. In: Vossen, P. (Editor). *EuroWordNet. A Multilingual Database with Lexical Semantic Networks*. Dordrecht: Kluwer Academic Publishers, pp 19-43.
- Andersson, Erik 1977. *Verbfrasens struktur i svenskan. En studie i aspekt, tempus, tidsadverbial och semantisk räckvidd*. Turku: Meddelanden från Stiftelsens för Åbo Akademi forskningsinstitut, 18.
- Agramann, Vered and Neal J. Pearlmutter 2002. Verb sense and verb subcategorization probabilities. In: Merlo, Paola and Suzanne Stevenson (Editors). *The Lexical Basis of Sentence Processing*. Amsterdam: John Benjamins, pp. 303-324.
- Arppe, Antti 2001. Focal points in frequency profiles - how some word forms in a paradigm are more significant than others in Finnish. Proceedings of the 6th Conference on Computational Lexicography and Corpus Research (COMPLEX), 28-30.6.2001, University of Birmingham, Birmingham, United Kingdom, pp. 1-7.
- Arppe, Antti 2002. The usage patterns and selectional preferences of synonyms in a morphologically rich language. In: Morin, Annie and Pascale Sébillot (Editors). *JADT-2002. 6th International Conference on Textual Data Statistical Analysis*, 13-15.3.2002, Vol. 1. Rennes: INRIA, pp. 21-32
- Arppe, Antti 2004. Every method makes a difference - describing the use of a Finnish synonym pair. Pre-proceedings of *International Conference on Linguistic Evidence*, January 29-31, 2004, Tübingen, Germany, pp. 137-138.
- Arppe, Antti 2005a. Morphological features as context in distinguishing semantically similar words. *Proceedings from the Corpus Linguistics Conference Series*, Vol. 1, no. 1, ISSN 1747-9398. Third Biennial Corpus Linguistics 2005 Conference, 14-17.7.2005, Birmingham, United Kingdom. Available online at URL: <http://www.corpus.bham.ac.uk/PCLC/>
- Arppe, Antti 2005b. The Very Long Way from Basic Linguistic Research to Commercially Successful Language Technology: the Case of Two-Level Morphology. In: *Inquiries into Words, Constraints, and Contexts. Festschrift for Kimmo Koskenniemi on his 60th Birthday*. Arppe, Antti, Lauri Carlson, Krister Lindén, Jussi Piitulainen, Mickael Suominen, Martti Vainio, Hanna Westerlund, and Anssi Yli-Jyrä (Editors). Stanford: CSLI Studies in Computational Linguistics ONLINE. Copestake, Ann (Series editor), pp. 2-17.
- Arppe, Antti 2006a. On the limits of generalizing from quantitative, corpus-based evidence in a morphologically rich language. Pre-proceedings of the *International Conference on Linguistic Evidence*. 2-4.2.2006, Tübingen, Germany, pp. 118-120. Available on-line at URL: <http://www.sfb441.uni-tuebingen.de/LingEvid2006/abstracts/arppe.pdf>

- Arppe, Antti 2006b. Frequency Considerations in Morphology, Revisited - Finnish Verbs Differ, Too. In: *A Man of Measure. Festschrift in Honour of Fred Karlsson in his 60th Birthday*. Suominen, Mickael, Antti Arppe, Anu Airola, Orvokki Heinämäki, Matti Miestamo, Urho Määttä, Jussi Niemi, Kari K. Pitkänen, and Kaius Sinnemäki (Editors). Special Supplement to *SKY Journal of Linguistics*, Volume 19/2006, pp. 175-189. Turku: Linguistic Association of Finland. Available online at URL: http://www.ling.helsinki.fi/sky/julkaisut/SKY2006_1/1.3.1.ARPPE.pdf
- Arppe, Antti 2006c. Complex phenomena deserve complex explanations. *Quantitative Investigations in Theoretical Linguistics* (QITL2) Conference, 1-2.6.2006, Osnabrück, Germany, pp. 8-11. Available on-line at URL: <http://www.cogsci.uni-osnabrueck.de/~qitl/>
- Arppe, Antti 2007. Multivariate methods in corpus-based lexicography. A study of synonymy in Finnish. In: Davies, Matthew, Paul Rayson, Susan Hunston, and Pernilla Danielsson (Editors). *Proceedings from the Corpus Linguistics Conference (CL2007)*, July 28-30, 2007, Birmingham, United Kingdom. Available on-line at: URL: <http://www.corpus.bham.ac.uk/corplingproceedings07/>
- Arppe, Antti 2008. Linguistic choices and probabilities – How much and what can linguistic theory explain? *Pre-Proceedings of the International Conference on Linguistic Evidence*, 31.1.-2.2.2008, Tübingen, Germany. Sonderforschungsbereich 441, University of Tübingen, pp. 9-12.
- Arppe, Antti, Mari Voipio, and Malene Würtz 2000. Creating Inflecting Electronic Dictionaries. Lindberg, Carl-Erik and Steffen Nordahl Lund (Editors). 17th Scandinavian Conference of Linguistics, Nyborg 20-22.8.1998. *Odense Working Papers in Language and Communication*, No 19, Vol. 1 (April 2000). Odense: University of Southern Denmark, pp. 1-11.
- Arppe, Antti and Järvikivi, Juhani 2002. Verbal Synonymy in Practice: Combining Corpus-Based and Psycholinguistic Evidence. Workshop on *Quantitative Investigations in Linguistics* (QITL1/QITL-2002), Osnabrück, Germany, 3-5.10.2002. Available on-line at URL: <http://www.cogsci.uni-osnabrueck.de/~qitl/>
- Arppe, Antti and Juhani Järvikivi 2007a. Take empiricism seriously! - In support of methodological diversity in linguistics [Commentary of Geoffrey Sampson (2007). Grammar without Grammaticality.] *Corpus Linguistics and Linguistic Theory*, Vol. 3, No. 1, pp. 99-109.
- Arppe, Antti and Juhani Järvikivi 2007b. Every method counts - Combining corpus-based and experimental evidence in the study of synonymy. *Corpus Linguistics and Linguistic Theory*. Vol. 3, No. 2, pp. 131-159.
- Atkins, Beryl T. S. 1987. Semantic ID tags: corpus evidence for dictionary senses. *Proceedings of the Third Annual Conference of the UW Centre for the New Oxford English Dictionary*, pp. 17-36.
- Atkins, Sue, Jeremy Clear, and Nicholas Ostler 1992. Corpus Design Criteria. *Literary and Linguistic Computing*, Vol. 7, No. 1, pp. 1-16.

- Atkins, Beryl T. S. and Beth Levin 1995. Building on a Corpus: A linguistic and lexicographical look at some near-synonyms. *International Journal of Lexicography*, 8:2, pp. 85–114.
- Baayen, R. Harald 2001. *Word Frequency Distributions*. Dordrecht: Kluwer Academic Publishers.
- Baayen, R. Harald 2007. *The languageR package*. Available on-line at URL: <http://cran.r-project.org/doc/packages/languageR.pdf>
- Baayen, R. Harald 2008. *Analyzing Linguistic Data. A Practical Introduction to Statistics*. Cambridge: Cambridge University Press.
- Baayen, R. Harald, Doug J. Davidson, and Douglas M. Bates (to appear 2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language* (Special issue on Emerging Data Analysis Techniques).
- Balota, David A. and James I. Chumbley 1985. The locus of word-frequency effects in the pronunciation task: Lexical access and/or production. *Journal of Memory and Language*, Vo. 24, pp. 89-106.
- Bard, Ellen Gurman, Dan Robertson, and Antonella Sorace 1996. Magnitude Estimation of Linguistic Acceptability. *Language*, Vol. 72, No. 1. (March 1996), pp. 32-68.
- Benjamini, Yoav and Daniel Yekutieli 2001. The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, Vol. 29, No. 4, pp.1165-1188.
- Biber, Douglas 1988. *Variation across speech and writing*. Cambridge: Cambridge University Press.
- Biber, Douglas 1993. Representativeness in Corpus Design. *Literary and Linguistic Computing*, Vol. 8, pp. 243-257. [Page number references from reprint in Sampson, Geoffrey and Diana McCarthy (Editors) 2005. *Corpus Linguistics. Readings in a Widening Discipline*. London/New York: Continuum, pp. 174-197]
- Biber, Douglas, Susan Conrad, and Randi Reppen 1998. *Corpus Linguistics: Investigating Language Structure and Use*. Cambridge: Cambridge University Press.
- Biber, Douglas and Jerry Kurjian 2007. Towards a taxonomy of web registers and text types: a multidimensional analysis. In: Hundt, Marianne, Nadja Nesselhauf, and Caroline Biewer (Editors). *Corpus Linguistics and the Web*. Language and Computers: Studies in Practical Linguistics, Amsterdam: Rodopi, pp. 109-131.
- Black, Jeremy A., Graham Cunningham, Jarle Ebeling, Esther Flückiger-Hawker, Eleanor Robson, Jon Taylor, Gábor Zólyomi 1998-2006. *The Electronic Text Corpus of Sumerian Literature*, Oxford. URL: <http://etcsl.orinst.ox.ac.uk/> (visited 29.11.2008).
- Bod, Rens, Jennifer Hay, and Stefanie Jannedy (Editors) 2003. *Probabilistic Linguistics*. Cambridge/London: MIT Press.
- Breiman, Leo 2001. Random Forests. *Machine Learning*, 45(1), pp. 5–32.

- Breiman, Leo and Adele Cutler 2005. Random Forests. URL: http://www.stat.berkeley.edu/users/breiman/RandomForests/cc_home.htm (visited 30.11.2008).
- Bresnan, Joan, Anna Cueni, Tatiana Nikitina, and R. Harald Baayen 2007. Predicting the Dative Alternation. In: *Cognitive Foundations of Interpretation*. Boume, G., I. Kraemer, and J. Zwarts. Amsterdam: Royal Netherlands Academy of Science, pp. 69-94.
- Bresnan, Joan 2006. Is syntactic knowledge probabilistic? Experiments with the English dative alternation. *Pre-proceedings of the International Conference on Linguistic Evidence*, 2-4 February 2006, Tübingen, Germany. Sonderforschungsbereich 441, University of Tübingen, pp. 3-10.
- Bresnan, Joan 2007. Is syntactic knowledge probabilistic? Experiments with the English dative alternation. In: Featherston, Sam and Wolfgang Sternefeld (Editors). *Roots: Linguistics in search of its evidential base*. Series: Studies in Generative Grammar. Berlin: Mouton de Gruyter.
- Buchanan, William 1974. Nominal and Ordinal Bivariate Statistics: The Practitioner's View. *American Journal of Political Science*, Vol. 18, No. 3. (August 1974), pp. 625-646.
- Bungarten, Theo 1979. Das Korpus als empirische Grundlage in der Linguistik und Literaturwissenschaft. In: Bergenholtz, Henning and Burkhard Schaefer (Editors). *Empirische Textwissenschaft. Aufbau und Auswertung von Text-Corpora*. Monografien Linguistik und Kommunikationswissenschaft 39. Königstein: Scriptor.
- Chafe, Wallace 1992. The importance of corpus linguistics to understanding the nature of language. In: Svartvik, Jan (Editor). *Directions in corpus linguistics*. Proceedings of the Nobel symposium 82, 4-8.8.1991, Stockholm. Trends in linguistics; Studies and monographs 65. Berlin: Mouton de Gruyter, pp. 79-97.
- Church, Kenneth, William Gale, Patrick Hanks, and Douglas Hindle 1991. Using Statistics in Lexical Analysis. In: Zernik, Uri (Ed.). *Lexical Acquisition: Exploiting On-Line Resources to Build a Lexicon*. Hillsdale: Lawrence Erlbaum Associates, pp. 115-164.
- Church, Kenneth, William Gale, Patrick Hanks, Douglas Hindle, Rosamund Moon 1994. Lexical substitutability. In: B. T. S. Atkind and Antonio Zampolli (Editors). *Computational Approaches to the Lexicon*. Oxford: Oxford University Press, pp. 153-177.
- Claridge, Claudia 2007. Constructing a corpus from the web: message boards. In: Hundt, Marianne, Nadja Nesselhauf, and Caroline Biewer. *Corpus Linguistics and the Web*. Language and Computers: Studies in Practical Linguistics, Amsterdam: Rodopi, pp. 87-108.
- Clear, Jeremy 1992. Corpus sampling. In: Leitner, Gerhard (Editor). *New Directions in English Language Corpora. Methodology, Results, Software Implementations*. Berlin/New York: Mouton de Gruyter, pp. 21-31.
- Clear, Jeremy, Gwyneth Fox, Gill Francis, Ramesh Krisnamurthy, and Rosamund Moon 1996. COBUILD: The State of the Art. *International Journal of Corpus Linguistics*, Vol. 1(2), pp. 303-314.

- Cochran, William G. 1952. The χ^2 Test of Goodness of Fit . *The Annals of Mathematical Statistics*, Vol. 23, No. 3 (September 1952), pp. 315-345.
- Cochran, William G. 1954. Some Methods for Strengthening the Common χ^2 Tests. *Biometrics*, Vol. 10, No. 4 (December 1954), pp. 417-451.
- Cohen, Jacob 1988. *Statistical power analysis for the behavioral sciences*, (2nd edition). Hillsdale: Lawrence Erlbaum Associates.
- Cohen, Jacob 1990. Things I Have Learned (So Far). *American Psychologist*, Vol. 45, No. 12, (December 1990), pp. 1304-1312.
- Cohen, Jacob 1992. A Power Primer. *Psychological Bulletin*, Vol. 112, No. 1, pp. 155-159.
- Cohen, Jacob 1994. The Earth is Round ($P < .05$). *American Psychologist*, Vol. 49, No. 12, pp. 997-1003.
- Cohen, Jacob, Patricia Cohen, Stephen G. West, and Leona S. Aiken 2003. *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences* (3rd edition). Mahwah: Lawrence Erlbaum Associates.
- Collot, Milena and Nancy Belmore 1996. Electronic language: A New Variety of English. In: *Computer-mediated Communication: linguistic, social and cross-cultural Perspectives*. Herring, Susan C. (Editor). Amsterdam: John Benjamins, pp. 13-28.
- Connexor 2007. List of morphological, surface-syntactic and functional syntactic features used in the linguistic analysis. [Web documentation] URL: <http://www.connexor.com/demo/doc/fifdg3-tags.html> (visited 29.5.2007) and URL: <http://www.connexor.com/demo/doc/enfdg3-tags.html> (visited 5.6.2007)
- Corbet, Michael and Michael Le Roy 2002. *Research Methods in Political Science: an Introduction Using Microcase*. Belmont: Wadsworth Publishing.
- Costner, Herbert L. 1965. Criteria for Measures of Association. *American Sociological Review*, Vol. 30, No. 3. (June 1965), pp. 341-353.
- Cox, D. R. 1958. The Regression Analysis of Binary Sequences. *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 20, No. 2, pp. 215-242.
- Cramér, Harald 1946. *Mathematical Methods in Statistics*. Princeton: Princeton University Press.
- Croft, William 2001. *Radical Construction Grammar. Syntactic Theory in Typological Perspective*. Oxford: Oxford University Press.
- Crombie, A. C. 1981. Philosophical Presuppositions and Shifting Interpretations of Galileo. In: Hintikka, Jaakko, David Gruender, and Evandro Agazzi. *Theory changes, Ancient Axiomatics, and Galileo's Methodology*. Proceedings of the 1978 Pisa Conference on the History and Philosophy of Science, Volume I. Dordrecht/Boston/London: Reidel, pp. 271-286.
- Cruse, D. Alan 1986. *Lexical Semantics*. Cambridge: Cambridge University Press.
- Cruse, D. Alan 2000. *Meaning in Language. An Introduction to Semantics and Pragmatics*. Oxford: Oxford University Press.

- Dahan, Daphne, James S. Magnusson, and Michael K. Tanenhaus 2001. Time course of frequency effects in spoken-word recognition: Evidence from eye-movements. *Cognitive Psychology*, Vol. 42, pp. 317-367.
- De Jong, Nivja 2002. *Morphological Families in the Mental Lexicon*. [PhD Dissertation] Nijmegen: MPI Series in Psycholinguistics.
- Divjak, Dagmar 2006. Ways on Intending. Delineating and Structuring Near-Synonyms. In: Gries, Stefan Th. and Anatol Stefanowitsch (Editors). *Corpora in cognitive linguistics*. Vol. 2: The syntax-lexis interface. Berlin: Mouton De Gruyter, pp. 19-56.
- Divjak, Dagmar and Stefan Th. Gries 2006. Ways of trying in Russian: Clustering and comparing behavioral profiles. *Corpus Linguistics and Linguistic Theory*. Vol 2(1), pp. 23-60.
- Divjak, Dagmar and Stefan Th. Gries. (forthcoming). Clusters in the mind? Converging evidence from Near-synonymy in Russian.
- Dowty, David 1991. Thematic Proto-Roles and Argument Selection. *Language*, Vol. 67, No. 3, pp. 547-619.
- Edmonds, Philip and Graeme Hirst 2002. Near-synonymy and Lexical Choice. *Computational Linguistics*, 28:2, pp. 105-144.
- Efron, Bradley 1979. Bootstrap Methods: Another Look at the Jackknife. *The Annals of Statistics*, Vol. 7, No. 1 (January 1979), pp. 1-26.
- Eliason, Scott A. 1993. *Maximum Likelihood Estimation. Logic and Practice*. Sage University Paper Series in Quantitative Applications in the Social Sciences 07-076. Newbury Park: Sage Publications.
- Ellis, Nick C. 2002. Frequency effects in language processing. A review with implications for theories of implicit and explicit language processing. *Studies in Second Language Acquisition*, Vol. 24, pp. 143-188.
- Estoup, J. B. 1916. *Gammes Sténographiques* (4th edition), Paris, France.
- Evert, Stefan and Marco Baroni 2006a. The zipfR library: Words and other rare events in R. *useR! 2006: The second R user conference*, Vienna, Austria. Available on-line at URL: <http://www.r-project.org/useR-2006/Slides/Evert+Baroni.pdf>
- Evert, Stefan and Marco Baroni 2006b. *The zipfR package*. Available on-line at URL: <http://cran.r-project.org/doc/packages/zipfR.pdf>
- Fagot, Anne M. 1981. Probabilities and Causes: On Life Tables, Causes of Death, and Etiological Diagnoses. In: Hintikka, Jaakko, David Gruender and Evandro Agazzi. *Probabilistic Thinking, Thermodynamics and the Interaction of the History and Philosophy of Science*. Proceedings of the 1978 Pisa Conference on the History and Philosophy of Science, Volume II. Dordrecht/Boston/London: Reidel, pp. 41-104.
- Featherston, Sam and Wolfgang Sternefeld 2007. *Roots: Linguistics in Search of its Evidential Base*. Berlin: Mouton de Gruyter.
- Featherston, Sam 2005. The Decathlon Model. In Kepser and Reis 2005a, pp. 187–208.

- Featherston, Sam 2007. Data in generative grammar: the stick and the carrot. *Theoretical Linguistics*, 33–3, pp. 269–318
- Fellbaum, Christiane 1998a. Introduction. In: Fellbaum, Christiane (Editor). *WordNet: An Electronic Lexical Database*. Cambridge: MIT Press, pp. 1-19.
- Fellbaum, Christiane 1998b. A Semantic Network of English Verbs. In: Fellbaum, Christiane (Editor). *WordNet: An Electronic Lexical Database*. Cambridge: MIT Press, pp. 69-104.
- Firth, J. R. 1957. A Synopsis of Linguistic Theory, 1930–1955. In: Firth, J. R. 1968. *Selected Papers of J. R. Firth 1952-1959*. London: Logmans, pp. 168-205.
- Flint, Aili 1980. *Semantic Structure in the Finnish Lexicon: Verbs of Possibility and Sufficiency*. Helsinki: Suomalaisen Kirjallisuuden Seura (SKST 360).
- Fortescue, Michael 2001. Thought about thought. *Cognitive Linguistics*, 12(1), pp. 15-45.
- Fox, John 1997. *Applied regression analysis, linear models, and related methods*. Thousand Oaks: Sage.
- Frank, Eibe and Stefan Kramer 2004. Ensembles of Nested Dichotomies for Multi-Class Problems. *Proceedings of the 21st International Conference on Machine Learning*, Banff, Canada.
- Fürnkranz, Johannes 2002. Round Robin Classification. *Journal of Machine Learning Research*, pp. 721-747.
- Garson, G. David 1975. *Handbook of Political Science Methods* (2nd edition). Boston: Holbrook Press.
- Garson, G. David 2007. *Statnotes: Topics in Multivariate Analysis*. URL: <http://www2.chass.ncsu.edu/garson/pa765/statnote.htm>. Visited Spring 2006 – Summer 2007.
- Goddard, Cliff 2002. The search for the shared semantic core of all languages. In: Goddard, Cliff and Anna Wierzbicka (Editors). *Meaning and Universal Grammar - Theory and Empirical Findings*. Volume I. Amsterdam: John Benjamins, pp. 5-40.
- Goddard, Cliff 2003. Thinking across languages and cultures: Six dimensions of variation. *Cognitive Linguistics*, 14-2/3, pp. 109-140.
- Goodman, Leo A. and William H. Kruskal 1954. Measures of Association for Cross-Classifications. *Journal of the American Statistical Association*, Vol. 49, No. 268 (December 1954), pp. 732-764.
- Gries, Stefan Th. 2002. Evidence in linguistics: Three approaches to genitives in English. In: Brend, Ruth M., William J. Sullivan and Arle R. Lommel (Editors). *LACUS Forum XXVIII: What Constitutes Evidence in Linguistics?* Fullerton: LACUS, pp. 17–31.
- Gries, Stefan Th. 2003a. *Multifactorial analysis in corpus linguistics: a study of particle placement*. London: Continuum.
- Gries, Stefan Th. 2003b. Towards a corpus-based identification of prototypical instances of constructions. *Annual Review of Cognitive Linguistics*, Vol. 1, pp. 1-27

- Gries, Stefan Th. 2003c. Testing the sub-test: a collocational-overlap analysis of English *-ic* and *-ical* adjectives. *International Journal of Corpus Linguistics*, Vol. 8(1), pp. 31-61.
- Gries, Stefan Th. 2006. Corpus-based methods and cognitive semantics: the many meanings of to run. In: Gries, Stefan Th. and Anatol Stefanowitsch (Editors). *Corpora in cognitive linguistics: corpus-based approaches to syntax and lexis*. TiLSM 172, Berlin/New York: Mouton de Gruyter, pp. 57-99.
- Gries, Stefan Th. 2007. Exploring variability within and between corpora: some methodological considerations. *Corpora*, Vol. 1 (2), pp. 109-51.
- Gries, Stefan Th. and Dagmar Divjak (forthcoming). Behavioral profiles: a corpus-based approach towards cognitive semantic analysis. In: Evans, Vyvyan and Stephanie S. Pourcel (Editors). *New directions in cognitive linguistics*. Amsterdam/Philadelphia: John Benjamins.
- Gries, Stefan Th., Beate Hampe and Doris Schönefeld 2005a. Converging evidence: Bringing together experimental and corpus data on the association of verbs and constructions. *Cognitive Linguistics*, 16-4, pp. 635-676.
- Gries, Stefan Th., Beate Hampe and Doris Schönefeld 2005b. Converging evidence II: More on the association of verbs and constructions. In Newman, John and Sally Rice (Editors). *Empirical and Experimental Methods in Cognitive/Functional Research*. Stanford: CSLI Publications.
- Gries, Stefan Th. and Anatol, Stefanowitsch 2004. Extending collostructional analysis. A corpus-based perspective on 'alternations'. *International Journal of Corpus Linguistics*, 9:1, pp. 97-129.
- Grondelaers, Stefan, Dirk Speelman, and Dirk Geeraerts 2002. Regressing on *er*. Statistical analysis of texts and language variation. In: Morin, Annie and Pascale Sébillot (Editors). *JADT-2002. 6th International Conference on Textual Data Statistical Analysis*, 13-15.3.2002, Vol. 1. Rennes: INRIA, pp. 335-346.
- Groth, Otto 1960. *Die unerkannte Kulturmacht. Grundlegung der Zeitungswissenschaft (Periodik)*, Vol. 1, Berlin: de Gruyter.
- Haarala, Risto and Lehtinen, Marja (Editors-in-Chief) 1990, 1993, 1994. *Suomen kielen perussanakirja (A-K), (L-R) and (S-Ö)*. Kotimaisten kielten tutkimuskeskus (KOTUS). Helsinki: Painatuskeskus.
- Haarala, Risto and Marja Lehtinen (Editors) 1997. *CD-Perussanakirja*. Kotimaisten kielten tutkimuskeskuksen julkaisu 94. Helsinki: Edita.
- Hacking, Ian 1981. From the Emergence of Probability to the Erosion of Determinism. In: Hintikka, Jaakko, David Gruender and Evandro Agazzi. *Probabilistic Thinking, Thermodynamics and the Interaction of the History and Philosophy of Science*. Proceedings of the 1978 Pisa Conference on the History and Philosophy of Science, Volume II. Dordrecht/Boston/London: Reidel, pp. 105-123.
- Hacking, Ian 1996. The Disunities of the Sciences. In: Galison, Peter and David J. Stump. *The Disunity of Science. Boundaries, Contexts, and Power*. Stanford: Stanford University Press, pp. 37-74.

- Hakulinen, Auli 2003. Ovatko puhuttu ja kirjoitettu kieli erkaantuneet toisistaan? *Kielikello* 1/2003.
- Hakulinen, Auli and Fred Karlsson 1979. *Nykysuomen lauseoppia* (3rd edition). Suomalaisen Kirjallisuuden Seura, Helsinki, Finland.
- Hakulinen, Auli, Anneli Kauppinen, Matti Leiwo, Heikki Paunonen, Anneli Räikkälä, Pauli Saukkonen, Valma Yli-Vakkuri, Jan-Ola Östman, and Irja Alho 1994. *Kieli ja sen kielioipit*. Helsinki: Opetusministeriö/Painatuskeskus.
- Hakulinen, Auli, Maria Vilkuna, Riitta Korhonen, Vesa Koivisto, Tarja-Riitta Heinonen, and Irja Alho 2004. *Iso suomen kielioippi*. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Hankala, Mari 2002. Sanomalehti mediakasvatuksessa. Kokeiluja perusopetuksen 7.-9. luokilla. Licentiate thesis in Journalism studies, University of Jyväskylä.
- Hanks, Patrick 1996. Contextual Dependency and Lexical Sets. *International Journal of Corpus Linguistics*, Vol. 1, No. 1, pp. 75-98.
- Harrell, Frank E. 2001. *Regression Modeling Strategies. With Applications to Linear Models, Logistic Regression and Survival Analysis*. New York: Springer-Verlag.
- Hartwig, Frederick and Brian A. Dearing 1979. *Exploratory Data Analysis*. Sage University Paper series on Quantitative Applications in the Social Sciences, 08-016. Beverly Hills/London: Sage Publications,.
- Heinonen, Tarja Riitta 2006. Kielioipin peruskäsitteitä: sanaluokat. *Kielikello* 3/2006, pp. 10-12.
- Herlin, Ilona 1997. *Suomen kielen koska-konjunktion merkitys ja merkityksenkehitys*. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Herlin, Ilona 1998. *Suomen 'kun'*. [PhD dissertation: Finnish 'kun']. SKST 712, Helsinki: Suomalaisen Kirjallisuuden Seura.
- Hochberg, Yosef 1988. A Sharper Bonferroni Procedure for Multiple Tests of Significance. *Biometrika*, Vol. 75, No. 4 (December 1988), pp. 800-802.
- Hoey, Michael 1991. *Patterns of Lexis in Text*. Oxford: Oxford University Press.
- Hoffman, Elaine B., Pranab K. Sen, and Clarice R. Weinberg 2001. Within-cluster resampling. *Biometrika*, Vol. 88, No. 4, pp. 1121-1134.
- Holland, Burt S. and Margaret DiPonzio Copenhaver 1988. Improved Bonferroni-Type Multiple Testing Procedures. *Psychological Bulletin*, Vol. 104, No. 1, pp. 145-149.
- Hommel, G. 1988. A stage-wise rejective multiple test procedure based in a modified Bonferroni test. *Biometrika*, Vol. 75, pp. 383-386.
- Hosmer, David W., Jr., and Stanley Lemeshow 1989. *Applied Regression Analysis*. New York: Wiley.
- Hosmer, David W., Jr., and Stanley Lemeshow 2000. *Applied Regression Analysis* (2nd edition). New York: Wiley.
- Howell, David C. 1999. *Fundamental Statistics for the Behavioral Sciences* (4th edition). Pacific Grove: Brooks/Cole Publishing Company.

- Huumo, Tuomas 1996. *Lokatiivit lauseen semanttisessa tulkinnassa. Ajan omistajan, paikan ja tilan adverbiaalien keskinäiset suhteet suomen kielessä*. [Locatives and Semantic Interpretation of the Sentence: On Mutual Relations of Adverbials Indicating Time, Possession, Space and (Internal) State in Finnish, PhD dissertation]. Publications of the Department of Finnish and General Linguistics 55, University of Turku.
- Inkpen, Diana 2004. *Building a Lexical Knowledge-Base of Near-Synonym Differences*. PhD dissertation, Department of Computer Science, University of Toronto.
- Inkpen, Diana and Graeme Hirst 2006. Building and Using a Lexical Knowledge-Base of Near-Synonym Differences. *Computational Linguistics* 32:2 (June 2006), pp. 223-262.
- Itkonen, Erkki and Ulla-Maija Kulonen (Editors-in-chief) 1992, 1995, 2000. *Suomen sanojen alkuperä. Etymologinen sanakirja 1-3*. Helsinki: Kotimaisten kielten tutkimuskeskus and Suomalaisen Kirjallisuuden Seura.
- Jantunen, Jarmo H. 2001. Tärkeä seikka ja keskeinen kysymys. Mitä korpuslingvistinen analyysi paljastaa lähisynonyyeistä? *Virittäjä* 105:2, pp. 170-192.
- Jantunen, Jarmo H. 2004. *Synonymia ja käännössuomi: korpusnäkökulma samamerkityksisyyden kontekstuaalisuuteen ja käännöskielen leksikaalisiin erityispiirteisiin*. [PhD Dissertation] University of Joensuu Publications in the Humanities 35, University of Joensuu.
- Jäppinen, Harri (Editor) 1989. *Synonymisanakirja*. [Nykysuomen sanakirja VII]. Werner Söderström, Porvoo, Finland.
- Jäppinen, Harri, Aarno Lehtola, E. Nelimarkka, and Matti Ylilammi 1983. Knowledge Engineering Approach to Morphological Analysis. *First Conference of the European Chapter of ACL*, Pisa, Italy, pp. 49-51.
- Jäppinen, Harri and Matti Ylilammi 1986. Associative Model of Morphological Analysis: An Empirical Inquiry. *Computational Linguistics*, Vol. 12, No 4, pp. 257-272.
- Järvinen, Timo and Pasi Tapanainen 1998. Towards and implementable dependency grammar. In: Kahane, Sylvain and Alain Polguère (Editors). *Proceedings of the Workshop on Processing of Dependency-based Grammars*, COLING-ACL'98, Montreal, Canada.
- Järvinen, Timo and Pasi Tapanainen 1997. *A Dependency Parser for English*. TR-1, Technical Reports of the Department of General Linguistics, University of Helsinki, Finland.
- Kangasniemi, Heikki 1992. *Modal Expressions in Finnish*. Studia Fennica, Linguistica 2. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Karlsson, Fred 1983. *Finnish Grammar*. Juva: Werner Söderström.
- Karlsson, Fred 1985. Paradigms and word forms. *Studia gramatyczne* VII, pp. 135–154.

- Karlsson, Fred 1986. Frequency considerations in morphology. *Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung* (ZPSK), Volume 39, Number 1, pp. 19–28.
- Karlsson, Fred 1990. Constraint grammar as a framework for parsing running text. In: Karlgren, Hans (Editor). *COLING -90: Papers Presented to the 13th International Conference on Computational Linguistics on the Occasion of the 25th Anniversary of COLING and the 350th Anniversary of Helsinki University*, Volume 3, pp. 168–173. Helsinki: Yliopistopaino.
- Karlsson, Fred 2008. *Finnish: An Essential Grammar* (2nd, thoroughly revised and expanded edition). London/New York: Routledge.
- Karlsson, Fred (to appear 2008). Early generative linguistics and empirical methodology. Kytö, Merja and Anke Lüdeling (Editors). *Handbook on Corpus Linguistics*. Berlin and New York: Mouton de Gruyter.
- Karlsson, Fred, Atro Voutilainen, Juha Heikkilä, and Arto Anttila (Editors) 1995. *Constraint Grammar: A language-independent system for parsing unrestricted text*. Natural Language Processing 4. New York: Mouton de Gruyter.
- Kaufman, Leonard and Peter J. Rousseeuw 1990. *Finding Groups in Data*. New York: John Wiley.
- Kenttä, Reetta 2004. From describing to explaining linguistic phenomena - combining methods. In: Kepser, Stephan and Marga Reis (Editors). *Pre-proceedings of the International Conference on Linguistic Evidence*, 29-31.1.2004, Tübingen, Germany. Tübingen: Sonderforschungsbereich 441 “Linguistic Data Structures”, pp. 157-158.
- Kepser, Stephan and Marga Reis (Editors) 2005a. *Linguistic Evidence. Empirical, Theoretical and Computational Perspectives*. Studies in Generative Grammar 85. Berlin/New York: Mouton de Gruyter.
- Kepser, Stephan and Marga Reis 2005b. Evidence in linguistics. In: Kepser, Stephan and Marga Reis (Editors) 2005a, pp. 1-6.
- Kilgariff, Adam and Gregory Grefenstette 2003. Introduction to the special issue on the web as the corpus. *Computational Linguistics* 29 (3), pp. 333-347.
- Klemola, Pertti 1981. *Helsingin Sanomat, sananvapauden monopoli*. Otava, Helsinki, Finland.
- Kohonen, Teuvo 1995. *Self-Organizing Maps*. Series in Information Sciences, Vol. 30. Heidelberg: Springer.
- Koskenniemi, Kimmo 1983. *Two-Level Morphology: A General Computational Model for Word-Form Recognition and Production*. [PhD dissertation]. Publications of the Department of General Linguistics, University of Helsinki, No. 11.
- Kotilainen, Lari 2007a. *Kiellon lumo. Kieltoisanaton kieltorakenne ja sen kiteytyminen*. Suomi 193. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Kotilainen, Lari 2007b. *Konstruktoiden dynamiikkaa*. [PhD dissertation].

- Kotilainen, Lari (forthcoming). Innovaatio, virhe vai vakiintunut konstruktio. In: Herlin, Ilona and Lari Kotilainen (Editors). *Verbit ja konstruktiot - tapaustutkimuksia suomesta*.
- Kukko, Mirjami 2003. Tekstiviestikeskustelu - kaaosta vai järjestystä? *Kielikello* 1/2003, pp. 11-13.
- Kvalseth, Tarald O. 1985. Cautionary Note about R^2 . *The American Statistician*, Vol. 39, No. 4, Part 1. (Nov., 1985), pp. 279-285.
- Kviz, Frederick J. 1981. Interpreting Proportional Reduction in Error Measures as Percentage of Variation Explained. *The Sociological Quarterly*, Vol. 22 (Summer 1981), pp. 413-420
- Lagus, Krista and Anu Airola 2001. Analysis of Functional Similarities of Finnish Verbs using the Self-Organizing Map. *Proceedings of the ESSLLI'2001*, Helsinki, Finland.
- Lebart, Ludovic, André Salem, Lisette Berry 1998. *Exploring Textual Data*. Text, Speech and Language Technology 4. Dordrecht: Kluwer Academic Publishers.
- Leech, Geoffrey 2005. Adding Linguistic Annotation. In: Wynne, Martin (Editor). *Developing Linguistic Corpora: a Guide to Good Practice*. Oxford: Oxbow Books, pp 17-29. Available online at <http://ahds.ac.uk/linguistic-corpora/> (Accessed 31.5.2007).
- Leech, Geoffrey 2007. New resources, or just better old ones? The Holy Grail of Representativeness. In: Hundt, Marianne, Nadja Nesselhauf, and Caroline Biewer. *Corpus Linguistics and the Web*. Language and Computers: Studies in Practical Linguistics, Amsterdam: Rodopi, pp. 133-149.
- Leech, Geoffrey 1993. Corpus Annotation Schemes. *Literary and Linguistic Computing*, Vol. 8, No. 4, pp. 275-281
- Lewin, Beverly A. and Yonatan Donner 2002. Communication in Internet message boards. *English Today* 71, Vol. 18, No. 3, pp. 29-37.
- Liebetrau, Albert M. 1983. *Measures of Association*. Sage University Paper series on Quantitative Applications in the Social Sciences, 07-032. Beverly Hills and London: Sage Publications.
- Lindén, Krister 2004. Evaluation of Linguistic Features for Word Sense Disambiguation with Self-Organized Maps. *Journal of Computers and the Humanities*, 38(4), pp. 417-435.
- Luostarinen, Heikki and Uskali, Turo 2004. *Suomalainen journalismi ja yhteiskunnan muutos 1980–2000*. Helsinki: Sitra, pp. 449-519. Available on-line at: URL: <http://www.sitra.fi/julkaisut/Heiskala.pdf> (Accessed 7.7.2007).
- Luukka, Minna-Riitta 2000. Sinulle on Postia! *Kielikello* 1/2000, pp. 24-28.
- Länsimäki, Maija 2007. Hiki voi virrata myös budjettiriihessä. *Helsingin Sanomat*, 4.3.2007. Available on-line at: URL: <http://www.kotus.fi/index.phtml?s=2198>.

- Löfberg, Laura, Dawn Archer, Scott Piao, Paul Rayson, Tony McEnery, Krista Varantola, and Jukka-Pekka Juntunen 2003. Porting an English semantic tagger to the Finnish language. In: Archer, Dawn, Paul Rayson, Andres Wilson, and Tony McEnery (Editors). *Proceedings of the Corpus Linguistics 2003 conference*. UCREL technical paper number 16. UCREL, Lancaster University, pp. 457-464.
- Makkonen-Craig, Henna 1996. *Yleispuhekielisyydet lehtikielessä*. [Unpublished Master's thesis]. Department of Finnish language and literature, University of Helsinki.
- Makkonen-Craig, Henna 2005. *Toimittajan läsnäolo sanomalehtitekstissä. Näkökulmia suomen kielen dialogisiin passiivilauseisiin*. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Malmberg, Tarmo 1984. *Journalismikritiikki*. Publications Series B (Reports), Vol. 15. Department of Journalism and Mass communication, University of Tampere, Finland.
- Manin, Dmitrii (submitted). Zipf's Law and Avoidance of Excessive Synonymy. [Manuscript submitted for publication]. *arXiv.org*. Available on-line at: <http://arxiv.org/abs/0710.0105>. Accessed 2.10.2007.
- Manning, Christopher D. and Hinrich Schütze 1999. *Foundations of statistical natural language processing*. Cambridge: MIT Press.
- Marcoccia, Michel 2004. On-line polylogues: conversation structure and participating framework in internet newsgroups. *Journal of Pragmatics* 36, pp. 115-145.
- Margolin, Barry H. and Richard J. Light 1974. An Analysis of Variance for Categorical Data II: Small Sample Comparisons with Chi Square and Other Competitors. *Journal of the American Statistical Association*, Vol. 69, No. 347 (September 1974), pp. 755-764.
- Mayerthaler, Willi 1981. *Morphologische Natürlichkeit* [Morphological Naturalness], Wiesbaden: Akademische Verlagsgesellschaft Athenaion.
- Menard, Scott 1995. Applied Logistic Regression Analysis. Sage University Paper Series on Quantitative Applications in the Social Sciences 07-106. Thousand Oaks: Sage Publications.
- Miller, George A. 1990. Nouns in WordNet: a lexical inheritance system. (revised August 1993). *International Journal of Lexicography*, 3 (4), pp. 245-264. Available on-line at: URL: <ftp://ftp.cogsci.princeton.edu/pub/wordnet/5papers.ps>
- Miller, George A. 1995. WordNet: A Lexical Database for English. *Communications of the ACM*, Vol. 38, No. 11 (November 1995), pp. 39-41.
- Miller, George A., Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J. Miller 1990. Introduction to WordNet: an on-line lexical database. *International Journal of Lexicography*, 3 (4), pp. 235-244. Available on-line at URL: <ftp://ftp.cogsci.princeton.edu/pub/wordnet/5papers.ps>
- Miller, Katherine J. 1998. Modifiers in WordNet. In: Fellbaum, Christiane (Editor). 1998, pp. 47-67.

- Mittlböck, M. and M. Schemper 2002. Explained Variation for Logistic Regression - Small Sample Adjustments, Confidence Intervals and Predictive Precision. *Biometrical Journal*, 44:3, pp. 263-272.
- Mittlböck, M. and M. Schemper 1996. Explained variation for logistic regression. *Statistics in Medicine*, 15. pp. 1987-1997.
- Mooney, Christopher Z. and Robert D. Duval 1993 *Bootstrapping: A Nonparametric Approach to Statistical Inference*. Sage University Paper Series on Quantitative Applications in the Social Sciences 07-095. Newbury Park: Sage Publications.
- Moran, Matthew D. 2003. Arguments for rejecting sequential Bonferroni in ecological studies. *OIKOS* 100:2, pp. 403-405.
- Mudraya, Olga, Bogdan Babych, Scott Piao, Paul Rayson, and Andrew Wilson 2006. Developing a Russian semantic tagger for automatic semantic annotation. In: *Proceedings of Corpus Linguistics 2006*, 10-14.10.2006, St. Petersburg, Russian Federation.
- Nakagawa, Shinichi 2004. A farewell to Bonferroni: the problems of low statistical power and publication bias. *Behavioral Ecology*, Vol. 15, No. 6, pp. 1044-1045.
- O'Keefe, Daniel J. 2003. Colloquy: Should Familywise Alpha be Adjusted? Against Familywise Alpha Adjustment. *Human Communication Research*, Vol. 9, No. 3 (July 2003), pp. 431-447.
- Olejnik, Stephen, Jianmin Li, Suchada Supattathum, and Carl J. Huberty 1997. Multiple Testing and Statistical Power With Modified Bonferroni Procedures. *Journal of Educational and Behavioral Statistics*, Vol. 22, No. 4 (Winter 1997), pp. 389-406.
- Pajunen, Anneli 1982. *Suomen kielen verbien leksikaalinen kuvaus*. Licentiate Thesis. Department of Finnish and Linguistics, University of Turku.
- Pajunen, Anneli 1988. *Verbien leksikaalinen kuvaus*. [Ph.D. dissertation]. Publications 18, Department of General Linguistics, University of Helsinki.
- Pajunen, Anneli 2001. *Argumenttirakenne: Asiantilojen luokitus ja verbien käyttäytyminen suomen kielessä*. Suomi 187. Helsinki: Suomalaisen Kirjallisuuden Seura.
- Palander, Marjatta 2005. Morfosyntaksi ja kielenhuolto - kirja- ja puhekielen rajankäyntiä. *Kielikello* 1/2005, pp. 14-17.
- Peduzzi, Peter, John Concato, Elizabet Kemper, Theodore R. Holford, and Alvan R. Feinstein 1996. A Simulation Study of the Number of Events per Variable in Logistic Regression Analysis. *Journal of Clinical Epidemiology*, Vol. 49, No. 12, pp. 1373-1379.
- Pedersen, Ted 1996. Fishing for Exactness. *Proceedings of the South-Central SAS Users Group Conference (SCSUG-96)*, 27-29.10.1996, Austin, Texas.
- Penke, Martina and Anette Rosenbach 2004a. Preface. *Studies in Language* 28:3, pp. 477-79.
- Penke, Martina and Anette Rosenbach 2004b. What counts as evidence in linguistics: An introduction. *Studies in Language* 28:3, pp. 480-526.

- Perneger, Thomas V. 1998. What's wrong with Bonferroni adjustments. *British Medical Journal*, 316, pp. 1236-1238.
- Pietilä, Kauko and Sondermann, Klaus 1994. *Sanomalehden yhteiskunta*. Tampere: Vastapaino.
- Pietilä, Kauko 1997. Society in the newspaper. In: Hyvärinen, Matti and Kauko Pietilä (Editors). *The institutes we live by*. Research Institute for Social Sciences, University of Tampere, Publications 17, pp. 19-58.
- Pulkkinen, Paavo 1992. "Lehtimetsävahingot havupuita pienempiä" [Damages to deciduous forests smaller than conifer trees". *Virittäjä*, 96: 2-3, pp. 286-290.
- Päiviö, Pia 2007. *Suomen kielen asti ja saakka. Terminatiivisten partikkelien synonymia, merkitys, käyttö, ja kehitys sekä asema kieliopissa* [PhD Dissertation]. Turun yliopiston suomalaisen ja yleisen kielitieteen laitoksen julkaisuja [Publications of the Department of Finnish and General Linguistics] 75, University of Turku.
- R Development Core Team 2007. *R: A language and environment for statistical computing*. URL: <http://www.R-project.org>.
- Rantanen, Leena 1999. *Ajan adverbiaalien referenssi ja definiittisyys*. [Unpublished Master's thesis]. Department of Finnish language and literature, University of Helsinki.
- Repo, Päivi 2003. Huipputeknologian Suomi puhuu maalaisvertauksin [High-tech Finland speaks with rural metaphors]. *Helsingin Sanomat*, 31.8.2003, p. A8.
- Reynolds, H. T. 1977. *Analysis of Nominal Data*. Sage University Paper series on Quantitative Applications in the Social Sciences, 08-007, Beverly Hills/London: Sage Publications.
- Rifkin, Ryan and Aldebaro Krakatau 2004. In Defense of One-Vs-All Classification. *Journal of Machine Learning Research*, pp. 101-141.
- Rodríguez, Horacio, Salvador Climent, Piek Vossen, Laura Bloksma, Wim Peters, Antonietta Alonge, Francesca Bertagna, and Adriana Roventini 1998. The Top-Down Strategy for Building EuroWordNet: Vocabulary Coverage, Base Concepts and Top Ontology. *Computers and the Humanities*, 32, pp. 117-132.
- Rom, Dror M. 1990. A Sequentially Rejective Test Procedure Based on a Modified Bonferroni Inequality. *Biometrika*, Vol. 77, No. 3 (September 1990), pp. 663-665.
- Roscoe, John T. and Jackson A. Byars 1971. An Investigation of the Restraints with Respect to Sample Size Commonly Imposed on the Use of the Chi-Square Statistic. *Journal of the American Statistical Association*, Vol. 66, No. 336 (December 1971), pp. 755-759.
- Rosenthal, Ralph A. and Robert Rosnow 1989. Statistical Procedures and the Justification of Knowledge in Psychological Science. *American Psychologist*, Vol 44, No. 10, pp. 1276-1284.
- Sadeniemi, Matti (Editor-in-chief) 1976 [1951-1961]. *Nykysuomen sanakirja* (5th edition). Porvoo: Suomalaisen Kirjallisuuden Seura/Werner Söderström.
- Sampson, Geoffrey 1995. *English for the computer: the SUSANNE corpus and analytic scheme*. Oxford: Clarendon Press.

- Sampson, Geoffrey 2007. Grammar without grammaticality. *Corpus Linguistics and Linguistic Theory*, 3:1 (Special issue: Grammar without grammaticality), pp. 1-32.
- Saukkonen, Pauli 2001. *Maailman hahmottaminen teksteinä. Tekstirakenteen ja tekstilajien teoriaa ja analyysia*. Helsinki: Helsinki University Press.
- Saukkonen, Pauli, Marjatta Haipus, Antero Niemikorpi, and Helena Sulkala 1979. *Suomen kielen taajuussanasto* [A Frequency Dictionary of Finnish]. Porvoo: Werner Söderström Osakeyhtiö.
- Sauri, Tuomo 2006a. Joukkoviestinnän talous ja kulutus. In: Sauri, Tuomo and Rauli Kohvakka (Editors). *Joukkoviestimet 2006* [Finnish Mass media]. Statistics Finland, Helsinki, Finland, pp. 117-142. Tables available on-line at: URL: http://www.stat.fi/til/jvie/2004/index_en.html (Accessed 7.7.2007).
- Sauri, Tuomo 2006b. Sanomalehdet. In: Sauri, Tuomo and Rauli Kohvakka (Editors). *Joukkoviestimet 2006* [Finnish Mass media]. Helsinki: Statistics Finland, pp. 271-296. Tables available on-line at: URL: http://www.stat.fi/til/jvie/2005/index_en.html (Accessed 7.7.2007).
- Sauri, Tuomo 2006c. Reading habits. In: Liikkanen, Mirja, Riitta Hanifi, and Ulla Hannula (Editors). *Individual choices, permanency of cultures. Changes in Leisure 1981-2002*. Helsinki: Statistics Finland, pp. 35-49.
- Sauri, Tuomo 2007. Sanomalehdet pystyvät vastaamaan ajan haasteisiin. *Hyvinvointikatsaus*, 2. Statistics Finland, pp. 28-31.
- Schreuder, Robert and R. Harald Baayen. 1997. How complex simplex words can be. *Journal of Memory and Language*, Vol. 37, pp. 118-119.
- SFNET co-ordinators 2007a. *Homepage of the sfnet newsgroups*. URL: <http://www.cs.tut.fi/sfnet/english.shtml> (Visited 7.7.2007).
- SFNET co-ordinators 2007b. *Sfnet-ryhmien ryhmäluettelo*. URL: <http://www.cs.tut.fi/sfnet/ryhmakuvaukset.shtml> (Visited 7.7.2007).
- Sinclair, John 1991. *Corpus, Concordance, Collocation*. Oxford: Oxford University Press.
- Sinclair, John (founding editor-in-chief) 2001. *Collins COBUILD English Dictionary for Advanced Learners* (3rd edition). Glasgow: HarperCollins.
- Sinclair, John 2005. Corpus and Text – Basic Principles. In: Wynne, Martin (Editor). *Developing Linguistic Corpora: a Guide to Good Practice*. Oxford: Oxbow Books, pp. 1-16. Available on-line at <http://ahds.ac.uk/linguistic-corpora/> [Visited 26.6.2007].
- Stassen, Leon 1985. *Comparison and Universal Grammar*. Oxford: Blackwell
- Stassen, Leon 2005. Comparative Constructions. In: Martin Haspelmath, Matthew Dryer, David Gil and Bernard Comrie (Editors). *The World Atlas of Language Structures*. Oxford: Oxford University Press, pp. 490-493. Available on-line at URL: <http://wals.info/>
- Stevens, S. S. 1946. On the theory of scales of measurement. *Science*, 103, pp. 677-680.

- Stevens, S. S. 1951. Mathematics, measurement and psychophysics. In S. S. Stevens (Editor), *Handbook of experimental psychology*. New York: Wiley, pp 1-49.
- Stubbs, Michael 1996 *.Text and Corpus Analysis. Computer-Assisted Studies of Language and Culture*. Oxford: Blackwell.
- Sulkala, Helena 1981. *Suomen kielen ajan adverbien semantiikkaa*. [PhD dissertation: Semantics of Finnish Temporal Adverbs]. Acta Universitatis Ouluensis, Series B, Humaniora 8, Philologica 3.
- Tapanainen, Pasi and Timo Järvinen 1997. A non-projective dependency parser. In: *Proceedings of the 5th Conference on Applied Natural Language Processing*, April 1997, Washington, D.C., Association of Computational Linguistics, pp. 64-71.
- Tesnière, Lucien 1959. *Éléments de syntaxe structurale*. Paris: Editions Klincksieck.
- Theil, Henri 1970. On the Estimation of Relationships Involving Qualitative Variables. *The American Journal of Sociology*, Vol. 76, No. 1 (July 1970), pp. 103-154.
- Thompson, Geoff 1998. Resonance in text. In: Sánchez-Macarro, Antonia and Roland Carter (Editors). *Linguistic choice across genres: Variation in spoken and written English*. Amsterdam: John Benjamins, pp. 29-46.
- Tinney, Steve and Michael Everson 2004-2007. *CompositeCuneiform* [font]. URL: <http://cdl.museum.upenn.edu/fonts.html> (visited 1.12.2008).
- Tolkien, J. R. R. 2007. *Húrinin lasten tarina*. Translated into Finnish by Kersti Juva from the English original *The Children of Húrin* (edited by Christoffer Tolkien). Helsinki: Werner Söderström Osakeyhtiö.
- Toivonen, Y. H., Erkki Itkonen, Aulis J. Joki, and Reino Peltola 1955, 1958, 1962, 1969, 1975, 1978, 1981. *Suomen Kielen Etymologinen Sanakirja I-VII*. Helsinki: Suomalais-Ugrilainen Seura.
- Valkonen, K., Harri Jäppinen, and Aarno Lehtola 1987. Blackboard-based dependency parsing. In: *Proceedings of IJCAI'87, Tenth International Joint Conference on Artificial Intelligence*, pp. 700-702.
- Vanhatalo, Ulla 2005. *Kyselytestit synonymian selvittämisessä – Sanastotietoutta kielenpuhujilta sähköiseen sanakirjaan*. [PhD. Dissertation]. Department of Fenno-Ugric Studies, University of Helsinki. Available on-line: URL: <http://ethesis.helsinki.fi/julkaisut/hum/suoma/vk/vanhatalo/kyselyte.pdf>
- Vanhatalo, Ulla 2003. Kyselytestit vs. korpuslingvistiikka lähisynonyymien semanttisten sisältöjen arvioinnissa – Mitä vielä keskeisestä ja tärkeästä?. *Virittäjä* 107:3, pp. 351-369.
- Váradi, Tamás 2001. The linguistic relevance of corpus linguistics. In: Rayson, Paul, Andres Wilson, Tony McEnery, A. Hardie, A. and S. Khoja (Editors). *Proceedings of the Corpus Linguistics 2001 Conference*. UCREL Technical Papers 13, Lancaster University, pp. 587-593.
- Velleman, Paul F. and Leland Wilkinson 1993. Nominal, Ordinal, Interval, and Ratio Typologies Are Misleading. *The American Statistician*, Vol. 47, No. 1 (Feb., 1993), pp. 65-72.

- Venables, William N. and Brian D. Ripley 1994. *Modern applied to statistics with S-PLUS* (Corrected fourth printing 1996). Statistics and Computing 2. New York/Berlin: Springer-Verlag.
- Vider, Kadri, Leho Paldre, Heili Orav, and Haldur Õim 1999. *The Estonian Wordnet*. EuroWordNet (LE-8328) Deliverable 2D014, Part B3. Available on-line at: URL: <http://www.ilc.uva.nl/EuroWordNet/docs.html> (visited 17.6.2007).
- Vossen, Piek (Editor) 1998a. *EuroWordNet: A Multilingual Database with Lexical Semantic Networks*. Dordrecht: Kluwer Academic Publishers.
- Vossen, Piek 1998b. Introduction to EuroWordNet. *Computers and the Humanities* 32, pp. 73-89.
- Weisberg, Herbert F. 1974. Models of Statistical Relationship. *The American Political Science Review*, Vol. 68, No. 4. (December 1974), pp. 1638-1655.
- Whorf, Benjamin 1956. *Language, Thought and Reality*. Cambridge: MIT Press. [Science and Linguistics (1940), s. 207–219; The relation of habitual thought and behavior to language (1939), s. 134–159.]
- Wikipedia contributors 2007a. Newsgroup. *Wikipedia, The Free Encyclopedia*. URL: <http://en.wikipedia.org/w/index.php?title=Newsgroup&oldid=142849690> (Visited 7.7.2007)
- Wikipedia contributors 2007b. Usenet. *Wikipedia, The Free Encyclopedia*. URL: <http://en.wikipedia.org/w/index.php?title=Usenet&oldid=145193135> (Visited 19.7.2007)
- Woods, Anthony, Paul Fletcher, and Arthur Hughes 1986. *Statistics in Language Studies*. Cambridge: Cambridge University Press.
- Yang, Charles 2008. The great number crunch. *Journal of Linguistics*, 44, pp. 205-228
- Ylikoski, Jussi 2005. Puhekielen morfologisten ja semanttisten innovaatioiden tutkimusnäkyimiä - esimerkkinä suomen tekeen- ja tekees-tyyppiset verbimuodot. *Puhe ja kieli*, 25:4, pp. 187–209. Available on-line at: URL: <http://cc oulu.fi/~jylikosk/puhekielen-vedos.pdf> (Visited 22.7.2007).
- Ylikoski, Jussi 2003. Havaintoja suomen ns. viidennen infinitiivin käytöstä. [With a summary in English: Remarks on the use of the proximative verb form (the so-called 5th infinitive) in Finnish]. *Sananjalka* 45, pp. 7-44. Available on-line at: URL: <http://cc oulu.fi/~jylikosk/5inf.pdf> (visited 31.5.2007).
- Zgusta, Ladislav 1971. *Manual of Lexicography*. Janua Linguarum, Series Maior 39. Prague: Academia.
- Zipf, George K. 1935. *The Psychobiology of Language*. Boston: Houghton Mifflin.
- Zipf, George K. 1949. *Human Behavior and the Principle of Least Effort*. Cambridge (Massachusetts): Addison-Wesley.